

**PONTIFICIA UNIVERSIDAD CATÓLICA DEL PERÚ**

**ESCUELA DE POSGRADO**



**Extensión al modelo DINA reparametrizado con covariable**

**TESIS PARA OPTAR POR EL GRADO DE MAGÍSTER EN  
ESTADÍSTICA**

**Presentado por:**

**Ing. Miguel Angel Sáenz Egúsquiza**

**Asesor: Dr. Luis Hilmar Valdivieso Serrano**

**Miembros del jurado:**

**Dr. Luis Hilmar Valdivieso Serrano**

**Dr. Cristian Luis Bayes Rodriguez**

**Mag. Enver Tarazona Vargas**

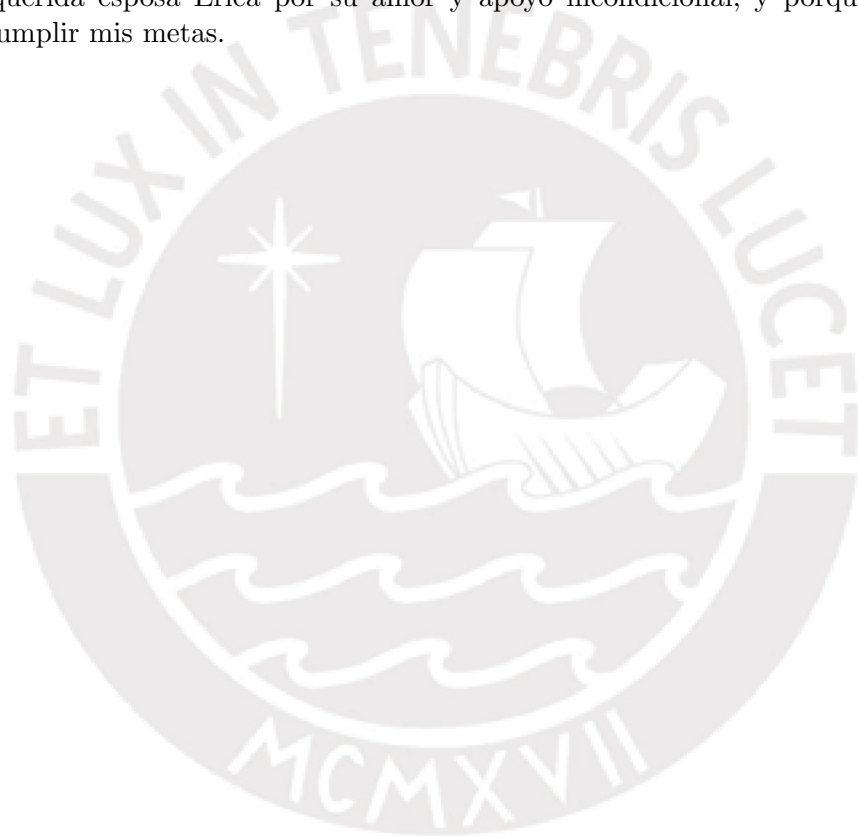
**Lima, Setiembre, 2020**

# Dedicatoria

Esta tesis está dedicado a mis padres Delia y Victoriano que desde el cielo me iluminan y siempre me alentaron a superarme.

A mi querido hijo Miguel Abdul que es mi motor y motivación para seguir adelante.

A mi querida esposa Erica por su amor y apoyo incondicional, y porque siempre me alienta a cumplir mis metas.





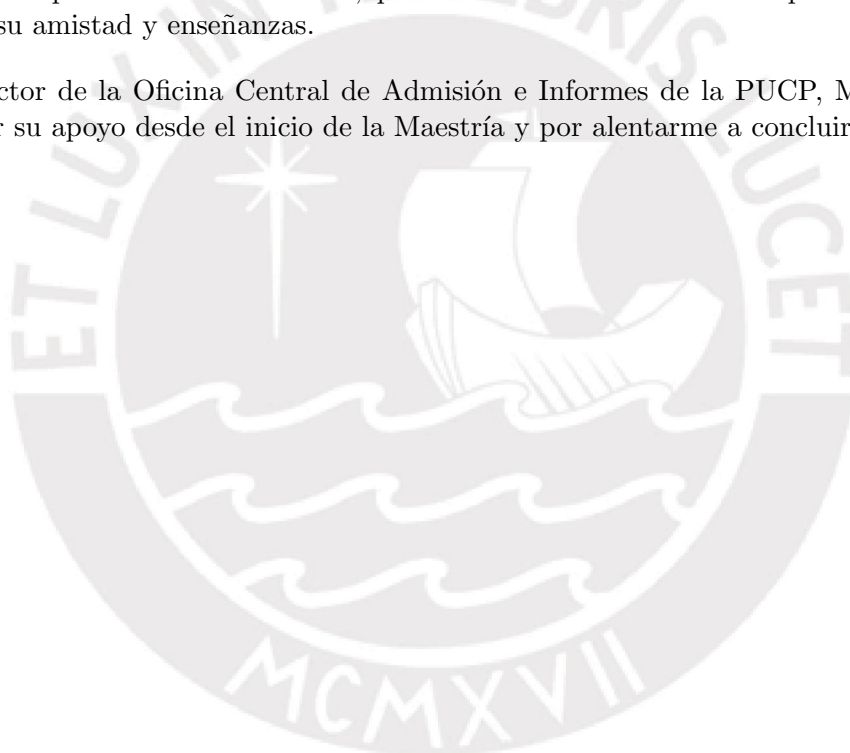
# Agradecimientos

Agradezco en primer lugar a mi asesor el Dr. Luis Valdivieso por su paciencia, por su apoyo constante, por sus enseñanzas y por el tiempo brindado durante todo el proceso de esta tesis.

A los docentes de la maestría en Estadística, por los conocimientos y buena disposición durante todo el programa que me motivaron a continuar en la investigación.

A mis compañeros de la Maestría, profesionales de diferentes disciplinas, quienes me brindaron su amistad y enseñanzas.

Al director de la Oficina Central de Admisión e Informes de la PUCP, Mag. Francisco Rivera, por su apoyo desde el inicio de la Maestría y por alentarme a concluirlo.



# Resumen

En el campo educacional, cuando los estudiantes resuelven problemas su habilidad en un tema particular puede influir en el desempeño de los mismos en un área de estudio similar pero diferente. Por ejemplo, la habilidad en ciencias podría tener un efecto en su dominio sobre las matemáticas, lo que a su vez afectará la forma en que los evaluados responden a las preguntas o ítems sobre matemáticas de una prueba. Por tanto, resulta natural examinar la relación entre el rendimiento en un área particular de estudio y el dominio de los atributos en un tema relacionado. Los modelos de diagnóstico cognitivo (CDM) proporcionan un marco ideal para realizar un análisis de este tipo, ya que clasifican a los examinados en perfiles de atributos que indican su dominio en las habilidades delimitadas permitiendo obtener información más específica con respecto a sus fortalezas y debilidades. Los CDM resuelven varias limitaciones de los métodos clásicos y los modelos de teoría de respuesta a ítems unidimensionales (TRI).

Para este estudio se amplía el marco de DINA al incorporar una covariable en un modelo de DINA reparametrizado. La covariable se puede especificar en dos niveles: en el nivel inferior, afectando la forma en que los evaluados resuelven los ítems (es decir, la probabilidad de respuesta), y en el nivel superior, influenciando en el dominio de los atributos (es decir, la clasificación latente). En esta tesis, se desarrolla teóricamente el modelo indicado desde el enfoque clásico. Para la estimación desarrollaremos el método de máxima verosimilitud y el método de la moda a posteriori vía el algoritmo de Esperanza-Maximización (EM) y de Newton-Raphson. Para tal fin, se realiza 4 estudios de simulación con la finalidad de observar en primer lugar el efecto de la covariable cuando afecta simultáneamente a los ítems y a los atributos, luego cuando la covariable afecta por separado a ambos, y también cuando la covariable no los afecta. Finalmente, se muestra su aplicación en la evaluación de la prueba de admisión a una Universidad.

**Palabras-clave:** Modelo DINA, Modelos de diagnóstico cognitivo (MDC), Modelo de regresión de clases latentes, estimación de parámetros.

# Índice general

<b>Lista de abreviaturas</b>	<b>VII</b>
<b>Índice de tablas</b>	<b>VIII</b>
<b>1. Introducción</b>	<b>1</b>
1.1. Planteamiento y justificación del tema . . . . .	1
1.2. Objetivos . . . . .	3
1.3. Organización del trabajo . . . . .	3
<b>2. Modelos de Clases Latentes y Diagnóstico Cognitivo</b>	<b>4</b>
2.1. Modelos de variables latentes . . . . .	5
2.2. Análisis de clases latentes . . . . .	6
2.3. El modelo de regresión de clases latentes . . . . .	7
2.4. El modelo regresión de clases latentes con predictores . . . . .	8
2.4.1. Estructura de la probabilidad y predictores lineales . . . . .	9
2.5. El modelo DINA . . . . .	10
2.6. El modelo RDINA . . . . .	11
<b>3. Extensión al modelo DINA con covariable</b>	<b>13</b>
3.1. El modelo RDINA con covariable . . . . .	13
<b>4. Estimación del modelo</b>	<b>17</b>
4.1. Función de verosimilitud . . . . .	17
4.2. Identificabilidad . . . . .	18
4.3. Estimación . . . . .	19
4.3.1. Estimación de parámetros por el método de la moda a posteriori . . . . .	19
4.3.2. El algoritmo EM . . . . .	19
4.3.3. Estimación de la varianza . . . . .	24
4.4. Selección del modelo . . . . .	25
<b>5. Estudio de Simulación</b>	<b>26</b>
5.1. Algoritmo para simular datos . . . . .	27
5.2. Criterios para evaluar la simulación . . . . .	27
5.3. Simulaciones . . . . .	28
5.3.1. Modelo RDINA con covariable afectando ítems y atributos . . . . .	28
5.3.2. Modelo RDINA con covariable afectando solo ítems . . . . .	31
5.3.3. Modelo RDINA con covariable afectando solo atributos . . . . .	33
5.3.4. Modelo RDINA sin covariable . . . . .	35

<b>6. Aplicación</b>	<b>37</b>
6.1. Estimación de parámetros del modelo RDINA con covariable afectando ítems y atributos . . . . .	40
6.2. Estimación de parámetros del modelo RDINA con covariable afectando solo ítems . . . . .	41
6.3. Estimación de parámetros del modelo RDINA con covariable afectando solo atributos . . . . .	42
6.4. Estimación de parámetros del modelo RDINA sin covariable . . . . .	43
6.5. Comentarios generales . . . . .	43
<b>7. Conclusiones y sugerencias</b>	<b>48</b>
7.1. Conclusiones . . . . .	48
7.2. Sugerencias para investigaciones futuras . . . . .	50
<b>A. Algoritmo Esperanza-Maximización (EM)</b>	<b>52</b>
A.1. Cálculo del paso Maximización(M) . . . . .	52
<b>B. Estudio de Simulación</b>	<b>55</b>
B.1. Valores iniciales, N=500, K=7, J=25 . . . . .	55
B.2. Matriz Q de habilidades . . . . .	56
B.3. Código R - Generación de Bases de Datos, N=500, K=7, J=25 . . . . .	56
B.3.1. Modelo RDINA con covariable afectando ítems y atributos . . . . .	56
B.3.2. Modelo RDINA con covariable afectando solo ítems . . . . .	57
B.3.3. Modelo RDINA con covariable afectando solo atributos . . . . .	57
B.3.4. Modelo RDINA sin covariable . . . . .	58
B.4. Código Latent Gold, con K=7, J=25 . . . . .	58
B.4.1. Modelo RDINA con covariable afectando ítems y atributos . . . . .	59
B.4.2. Modelo RDINA con covariable afectando solo ítems . . . . .	59
B.4.3. Modelo RDINA con covariable afectando solo atributos . . . . .	60
B.4.4. Modelo RDINA sin covariable . . . . .	60
<b>C. Aplicación - Códigos Latent Gold</b>	<b>62</b>
C.1. Código Latent Gold, con N=727, K=8, J=36 . . . . .	62
C.1.1. Modelo RDINA con covariable afectando ítems y atributos . . . . .	62
C.1.2. Modelo RDINA con covariable afectando solo ítems . . . . .	63
C.1.3. Modelo RDINA con covariable afectando solo atributos . . . . .	64
C.1.4. Modelo RDINA sin covariable . . . . .	65
<b>D. Algunas preguntas en la aplicación del examen de admisión</b>	<b>66</b>
D.1. Competencia: Redacción . . . . .	66
D.2. Competencia: Lectura . . . . .	69
<b>Bibliografía</b>	<b>72</b>

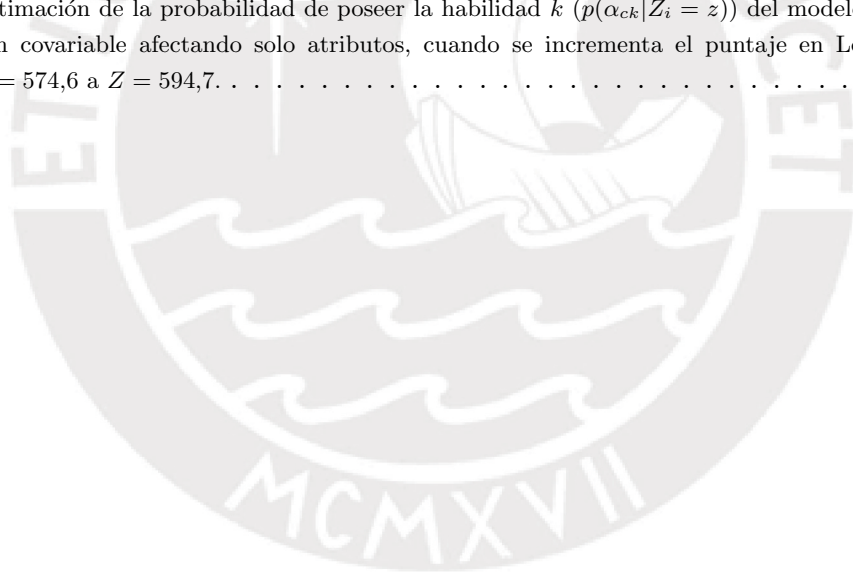
# Lista de abreviaturas

CDM	<i>Cognitive Diagnosis Model.</i>
MDC	Modelos de Diagnóstico Cognitivo.
DINA	<i>Deterministic Input Noisy AND gate.</i>
RDINA	<i>Reparameterized Deterministic Input Noisy AND gate.</i>
EM	Esperanza - Maximización.
LCA	<i>Latent Class Analysis.</i>
LCR	<i>Latent Class Regression.</i>
LCRP	<i>Latent Class Regression with Predictors.</i>
MSE	<i>Mean Squared Error.</i>
AIC	Criterio de Información de Akaike.
BIC	Criterio de Información Bayesiano.
GLM	<i>Generalized linear model.</i>
MCMC	<i>Markov Chain Monte Carlo.</i>
MLE	<i>Maximum likelihood estimation.</i>
MV	<i>Máxima verosimilitud.</i>
MAP	Maximización a posteriori.
NR	Newton-Raphson.
TIMSS	<i>Trends in International Mathematics and Science Study.</i>

# Índice de tablas

5.1. Estimación del vector de parámetros <b>b</b> del modelo RDINA con covariable afectando ítems y atributos en la simulación . . . . .	28
5.2. Estimación del vector de parámetros <b>h</b> del modelo RDINA con covariable afectando ítems y atributos en la simulación . . . . .	29
5.3. Estimación del vector de parámetros <b>f</b> del modelo RDINA con covariable afectando ítems y atributos en la simulación . . . . .	29
5.4. Estimación del vector de parámetros <b>d</b> del modelo RDINA con covariable afectando ítems y atributos en la simulación . . . . .	30
5.5. Estimación del vector de parámetros $\ell$ del modelo RDINA con covariable afectando ítems y atributos en la simulación . . . . .	30
5.6. Estimación del vector de parámetros <b>b</b> del modelo RDINA con covariable afectando solo ítems en la simulación . . . . .	31
5.7. Estimación del vector de parámetros <b>f</b> del modelo RDINA con covariable afectando solo ítems en la simulación . . . . .	31
5.8. Estimación del vector de parámetros <b>d</b> del modelo RDINA con covariable afectando solo ítems en la simulación . . . . .	32
5.9. Estimación del vector de parámetros $\ell$ del modelo RDINA con covariable afectando solo ítems en la simulación . . . . .	33
5.10. Estimación del vector de parámetros <b>b</b> del modelo RDINA con covariable afectando solo atributos en la simulación . . . . .	33
5.11. Estimación del vector de parámetros <b>h</b> del modelo RDINA con covariable afectando solo atributos en la simulación . . . . .	34
5.12. Estimación del vector de parámetros <b>f</b> del modelo RDINA con covariable afectando solo atributos en la simulación . . . . .	34
5.13. Estimación del vector de parámetros <b>d</b> del modelo RDINA con covariable afectando solo atributos en la simulación . . . . .	35
5.14. Estimación del vector de parámetros <b>b</b> del modelo RDINA sin covariable en la simulación .	35
5.15. Estimación del vector de parámetros <b>f</b> del modelo RDINA sin covariable en la simulación .	36
5.16. Estimación del vector de parámetros <b>d</b> del modelo RDINA sin covariable en la simulación .	36
6.1. Competencias evaluadas en la prueba de admisión en la aplicación . . . . .	37
6.2. Habilidades definidas en la prueba de admisión en la competencia Redacción . . . . .	38
6.3. División de la prueba de admisión en la competencia Redacción por temas y subtemas . .	38
6.4. Número de parámetros a estimar según modelo . . . . .	39
6.5. Matriz Q de habilidades de la competencia Redacción . . . . .	39
6.6. Resultados generales de los modelos ajustados en la aplicación . . . . .	40
6.7. Estimación de los vectores de parámetros <b>b</b> y <b>h</b> del modelo RDINA con covariable afectando ítems y atributos en la aplicación . . . . .	40
6.8. Estimación de los vectores de parámetros <b>f</b> , <b>d</b> y <b>l</b> del modelo RDINA con covariable afectando ítems y atributos en la aplicación . . . . .	41

6.9. Estimación del vector de parámetros <b>b</b> del modelo RDINA con covariable afectando solo ítems en la aplicación . . . . .	42
6.10. Estimación de los vectores de parámetros <b>f</b> , <b>d</b> y <b>l</b> del modelo RDINA con covariable afectando solo ítems en la aplicación . . . . .	42
6.11. Estimación de los vectores de parámetros <b>b</b> y <b>h</b> del modelo RDINA con covariable afectando solo atributos en la aplicación . . . . .	43
6.12. Estimación de los vectores de parámetros <b>f</b> y <b>d</b> del modelo RDINA con covariable afectando solo atributos en la aplicación . . . . .	44
6.13. Estimación del vector de parámetros <b>b</b> del modelo RDINA sin covariable en la aplicación . . . . .	44
6.14. Estimación de los vectores de parámetros <b>f</b> y <b>d</b> del modelo RDINA sin covariable en la aplicación . . . . .	45
6.15. Estimación de los parámetros de adivinación ( $g$ ) y desliz ( $s$ ) para un $Z = 574,6$ , que equivale a 28 ítems respondidos correctamente en promedio por los examinados en la prueba de Lectura. . . . .	45
6.16. Estimación de los parámetros de adivinación ( $g$ ) y desliz ( $s$ ) del modelo RDINA con covariable afectando ítems y atributos, cuando se incrementa el puntaje en Lectura de $Z = 574,6$ a $Z = 594,7$ . . . . .	47
6.17. Estimación de los parámetros de adivinación ( $g$ ) y desliz ( $s$ ) del modelo RDINA con covariable afectando solo ítems, cuando se incrementa el puntaje en Lectura de $Z = 574,6$ a $Z = 594,7$ . . . . .	47
6.18. Estimación de la probabilidad de poseer la habilidad $k$ ( $p(\alpha_{ck} Z_i = z)$ ) del modelo RDINA con covariable afectando ítems y atributos, cuando se incrementa el puntaje en Lectura de $Z = 574,6$ a $Z = 594,7$ . . . . .	47
6.19. Estimación de la probabilidad de poseer la habilidad $k$ ( $p(\alpha_{ck} Z_i = z)$ ) del modelo RDINA con covariable afectando solo atributos, cuando se incrementa el puntaje en Lectura de $Z = 574,6$ a $Z = 594,7$ . . . . .	47





# Capítulo 1

## Introducción

### 1.1. Planteamiento y justificación del tema

En el campo educacional, la habilidad de un estudiante en un tema particular, por ejemplo en ciencias fácticas como la física y la química, podría pensarse que puede influir en su desempeño de un área de estudio similar pero diferente como las matemáticas. Es decir, la habilidad en ciencias podría tener un efecto en su dominio sobre las matemáticas, lo que a su vez afectará la forma en que los evaluados responden a las preguntas o ítems sobre matemáticas de una prueba. En particular, se considera que las matemáticas y la ciencia tienen una relación estructural y funcional, por lo que las matemáticas pueden usarse como una herramienta para la ciencia, y la ciencia también puede funcionar como un estímulo para futuros descubrimientos matemáticos (Li et al.; 2002).

Por tanto, resulta natural examinar la relación que puede tener la capacidad científica en las habilidades matemáticas, lo que a su vez también puede afectar la forma en que los evaluados responden a los ítems de matemáticas. La relación puede no ser necesariamente causal, pero la dirección y la magnitud de la asociación puede ser información valiosa. En efecto, la pregunta que puede ser valiosa para los investigadores es si las características de la ciencia están relacionadas con habilidades matemáticas específicas; a la inversa, también pueden preguntar si la habilidad matemática afecta atributos específicos en la ciencia. Este concepto se extiende más allá del cálculo de correlaciones simples de rendimiento de los estudiantes en matemáticas y ciencias que tradicionalmente se realizan en muchos estudios. Si la capacidad y el conocimiento de un área temática puede influir en el dominio de las habilidades y los conceptos en otra área, la identificación de estas habilidades específicas pueden mejorar no solo la enseñanza sino también la retroalimentación a los estudiantes que carecen de dichas habilidades. Por tanto, este estudio es una oportunidad para examinar el efecto que el conocimiento en cualquiera de las asignaturas afecta la probabilidad de un evaluado para dominar las habilidades o mejorar la resolución de preguntas en un tema diferente, con el que comparte ciertas competencias.

En la presente tesis se examina en la aplicación la relación que pueda tener la capacidad en Lectura de los examinados en las habilidades en Redacción evaluada en una prueba de admisión, lo cual puede afectar la forma en que los evaluados responden los ítems de Redacción.

Los modelos de diagnóstico cognitivo (CDMs) proporcionan un marco ideal para realizar un análisis de este tipo, ya que clasifican a los examinados en perfiles de habilidades que indican su dominio en dichas habilidades delimitadas. Para este estudio se busca incorporar a los CDM el uso de una covariable. Para simplificar se trabajará con uno de los modelos CDM, el modelo DINA (deterministic input noisy and gate), Junker y Sijtsma (2001), en el cual la covariable se podría especificar en dos niveles: en el nivel inferior, afectando la forma en que los evaluados resuelven los ítems (es decir, la probabilidad de respuesta), y en el nivel



superior, influenciando en el dominio de las habilidades (es decir, la clasificación latente).

El modelo DINA requiere la identificación de habilidades necesarias para responder un ítem lo cual requiere la construcción de una matriz  $Q$  (Tatsuoka; 1985). Sea  $Y_{ij}$  una variable aleatoria correspondiente a la respuesta del examinado  $i$ , ítem  $j$ , y sea  $\alpha = (\alpha_1, \alpha_2, \dots, \alpha_K)'$  un vector de variables dicotómicas que indica la presencia o ausencia de  $K$  habilidades en estudio. La matriz  $Q$  es una matriz binaria de orden  $J$  por  $K$  que especifica la relación entre los ítems y las habilidades; en la cual un valor de 1 indica que una habilidad particular es requerida para responder el ítem, mientras que un valor de 0 representa que la habilidad particular no es necesaria.

El modelo DINA es conjuntivo, pues asume que todas las habilidades especificadas son necesarias para que el examinado resuelva el ítem. Esto se indica mediante la variable latente binaria  $\eta_{ij} = \prod_{k=1}^K \alpha_{ik}^{q_{jk}}$ , que clasifica si el examinado  $i$  domina todas las habilidades requeridos para el ítem  $j$ . El modelo DINA calcula la probabilidad de que el examinado  $i$  resuelva el ítem  $j$  correctamente dada su capacidad.

En este modelo se introducen dos parámetros: el desliz  $s_j$  y la adivinación  $g_j$  (Junker y Sijtsma; 2001). Los estudiantes que poseen todas las habilidades requeridas para un ítem pueden tener un desliz y responder incorrectamente al ítem, y los estudiantes que no poseen todas las habilidades requeridas para un ítem pueden adivinar y responder correctamente el ítem. El modelo DINA define el parámetro de desliz como  $s_j = p(Y_{ij} = 0 | \eta_{ij} = 1)$  y el parámetro de adivinación como  $g_j = p(Y_{ij} = 1 | \eta_{ij} = 0)$  para el ítem  $j$ .

DeCarlo (2011) ha propuesto una nueva reparametrización de este modelo dada por los parámetros  $f_j$  y  $d_j$  a través de:

$$\begin{aligned} g_j &= \exp(f_j) / (1 + \exp(f_j)) \\ s_j &= 1 - \exp(f_j + d_j) / (1 + \exp(f_j + d_j)) \end{aligned}$$

Al tratar el modelo DINA como un modelo de clase latente, la probabilidad de las respuesta del examinado  $i$  a los  $j$  ítems se puede modelar de la siguiente manera (DeCarlo; 2011):

$$p(y_{i1}, y_{i2}, \dots, y_{iJ}) = \sum_{c=1}^C p(\alpha_c) p(y_{i1}, y_{i2}, \dots, y_{iJ} | \alpha_c) = \sum_{c=1}^C p(\alpha_c) \prod_{j=1}^J p(y_{ij} | \alpha_c). \quad (1.1)$$

Cuando se aplica el supuesto restrictivo de independencia para la estructura de habilidades  $p(\alpha_c)$  se convierte en:

$$p(\alpha_c) = \prod_{k=1}^K p(\alpha_k) = \prod_{k=1}^K \exp(b_k) / (1 + \exp(b_k)). \quad (1.2)$$

Como se indica en (1.2), el modelo logístico para  $p(\alpha_k) = \exp(b_k) / (1 + \exp(b_k))$  utiliza un parámetro de dificultad de la habilidad  $b_k$ . Al introducir una covariable discreta o continua,  $Z$ , la probabilidad de respuesta en (1.1) se modifica de la siguiente manera:

$$p(y_{i1}, y_{i2}, \dots, y_{iJ} | Z) = \sum_{c=1}^C p(\alpha_c | Z) \prod_{j=1}^J p(y_{ij} | \alpha_c, Z). \quad (1.3)$$

En (1.3) se representa la probabilidad de respuesta condicionado a la covariable  $Z$ , que posteriormente se puede separar en dos términos con la finalidad de examinar el efecto de la covariable en las probabilidades de respuesta,  $p(y_{ij} | \alpha_c, Z)$ , o en la probabilidad de la habilidad,  $p(\alpha_c | Z)$ .

Cuando la covariable condiciona las probabilidades de respuesta, ésta puede cambiar la estimación de los parámetros de adivinación y desliz. Asimismo, cuando la covariable

condiciona las probabilidades de la habilidad, ésta se convierte en un predictor de los patrones de las habilidades que afectan a la membresía de la clase latente.

La presente tesis se enfocará en el desarrollo de este modelo desde el punto de vista clásico.

## 1.2. Objetivos

El objetivo general de esta investigación será el de estudiar un modelo de diagnóstico cognitivo RDINA covariable, y mostrar su aplicación en un conjunto de datos reales. De manera específica:

- Estudiar y presentar los fundamentos de los modelos de diagnóstico cognitivo y su extensión con el uso de una covariable.
- Estudiar y presentar el proceso de estimación del modelo haciendo énfasis en el enfoque clásico.
- Realizar estudios de simulación con la finalidad de recuperar los parámetros del modelo RDINA.
- Aplicar el modelo a una evaluación del examen de admisión a una universidad.

## 1.3. Organización del trabajo

En el capítulo 2, se presenta una revisión teórica de los modelos de clases latentes (LCA, LCR, LCRP), luego de los modelos de diagnóstico cognitivo, donde se presenta el modelo más popular denominado DINA.

En el capítulo 3, se presenta una ampliación en el marco del modelo DINA al incluir una covariable en base a su modelo reparametrizado RDINA, considerando sus propiedades y supuestos básicos de independencia.

En el capítulo 4, se presenta la estimación del modelo RDINA con covariable usando los métodos de máxima verosimilitud y de la moda a posteriori, vía el algoritmo de Esperanza-Maximización (EM) y de Newton-Raphson.

En el capítulo 5, se presenta un estudio de simulación considerando cuatro escenarios, con la finalidad de examinar la precisión en la recuperación de los parámetros.

En el capítulo 6, se presenta la aplicación del modelo a un conjunto de datos correspondiente a una prueba de admisión rendida en una universidad, con la finalidad de evaluar el ajuste de los modelos bajo los cuatro escenarios y analizar el efecto de la covariable.

Finalmente, en el capítulo 7, se presenta algunas conclusiones obtenidas en este trabajo, asimismo se plantea sugerencias para trabajos posteriores.

En el Anexo se presenta el algoritmo de Esperanza-Maximización (Apéndice A). Se incluye además los códigos para el estudio de simulación de los cuatro escenarios (Apéndice B), y también los códigos para la aplicación del modelo (Apéndice C). Asimismo, se muestra algunas preguntas del examen de admisión usadas en la aplicación (Apéndice D).

## Capítulo 2

# Modelos de Clases Latentes y Diagnóstico Cognitivo

En las ciencias sociales, del comportamiento y de la salud pueden presentarse diferentes subgrupos, tipos o categorías de individuos que presentan características disímiles entre sí. Un ejemplo es brindado por [Coffman et al. \(2007\)](#) quienes identificaron subgrupos de estudiantes del último año de secundaria de EE.UU. que tenían diferentes motivaciones para beber. Otro ejemplo, es dado por [Kessler et al. \(1998\)](#) que sobre la base de una muestra de residentes de EE. UU. entre 15 y 54 años que participaron en la Encuesta Nacional de Comorbilidad ([Kessler et al.; 1994](#)), identificaron dos tipos de fobias sociales. Cada uno de estos estudios utilizó un análisis de clases latentes (LCA) para identificar tales subgrupos a partir de los datos empíricos.

[Huang y Bandeen-Roche \(2004\)](#) extienden el modelo de clase latente permitiendo que tanto la distribución de la variable de clase subyacente como las distribuciones dentro de la clase de los indicadores medidos estén relacionados con variables independientes a nivel individual (la extensión de regresión del análisis de clase latente, LCRP). Esta idea no es nueva en sí misma. Se han desarrollado modelos de regresión bastante generales para describir la relación entre las covariables y la variable subyacente ([Dayton y Macready; 1988](#); [Van der Heijden et al.; 1996](#); [Bandeen-Roche et al.; 1997](#)), con la finalidad de medir los efectos de las variables independientes sobre la subyacente, o modelos para describir la relación entre covariables y los indicadores medidos ([Melton et al.; 1994](#)), con la finalidad de ajustar las características asociadas con la medición, evitando así posibles errores de clasificación de las categorías de las variables subyacentes.

En este capítulo, comenzaremos revisando conceptos acerca de los modelos LCA. Continuaremos luego revisando los modelos de diagnóstico cognitivo (MDC) que son modelos psicométricos, cuya finalidad es describir o diagnosticar el comportamiento de los examinados por medio de clases o perfiles latentes, obteniéndose información más específica acerca de sus fortalezas y debilidades, permitiendo una efectiva medición del aprendizaje y progreso del estudiante, apoyando el diseño de una mejor instrucción y posiblemente una intervención para hacer frente a las necesidades individuales y grupales ([Huebner; 2010](#)). Uno de los MDC más populares es el denominado DINA, que tuvo su primera aparición en [Haertel \(1989\)](#) enfocado básicamente en el campo educacional. Este modelo considera solo respuestas observadas dicotómicas de parte de los examinados teniendo como principal restricción que estos deben dominar necesariamente todas las habilidades requeridas para responder correctamente cada ítem. Además, el modelo permite estimar parámetros para los ítems, los cuales son denominados de adivinación y deslíz.

## 2.1. Modelos de variables latentes

El análisis de clases latentes (LCA) involucra el uso de variables latentes. Dichas variables no se miden directamente, sino en forma indirecta por medio de una o más variables observables. A diferencia de las variables latentes, las variables observables están sujetas a error. Las variables latentes a menudo se denominan constructos, particularmente en psicología y campos relacionados (Pedhazur y Pedhazur Schmelkin; 1991).

En el LCA, cada variable latente es categórica y está compuesta por un conjunto de modalidades que presentan una distribución multinomial. Las variables observables son una función de la variable latente y del error, y los cambios en ellas se asumen que son inducidos por cambios en los constructos latentes. Para un mejor entendimiento, se ilustra en la siguiente figura una variable latente hipotética la cual está representada por un óvalo. Las variables observables que miden la variable latente están representadas por cuadrados etiquetados como  $X_1$ ,  $X_2$  y  $X_3$ . Los círculos que contienen las letras  $e_1$ ,  $e_2$  y  $e_3$  representan los componentes del error asociados con  $X_1$ ,  $X_2$  y  $X_3$ , respectivamente. Hay flechas que se ejecutan desde la variable latente a cada variable observable, así como flechas que se dirigen desde cada componente de error a cada variable observable. Estas flechas representan un concepto importante que subyace a todos los modelos de variables latentes, incluido el LCA: las causas de las variables observables son las variables latentes y el error. Es importante notar, que el flujo causal es de la variable latente a la variable observable no al revés, es decir, las variables observables miden las variables latentes, pero las variables observables no causan las variables latentes.

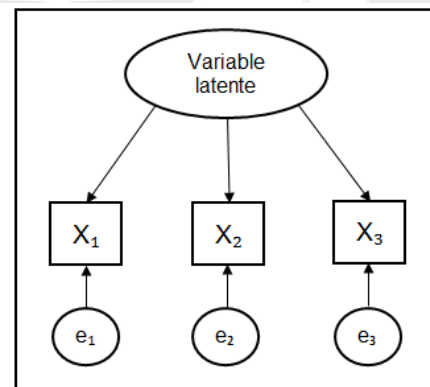


Figura 2.1: Variable latente con tres variables observables

El cuadro siguiente muestra una visión general de como el LCA se relaciona con algunos otros modelos de variables latentes. Como se aprecia, los modelos de variables latentes pueden organizarse según si la variable latente es categórica o continua, y si las variables observables se tratan como categóricas o continuas (Collins y Lanza; 2009).

	<b>Variables latentes continuas</b>	<b>Variables latentes categóricas</b>
<b>Variables observables continuas</b>	Análisis factorial	Análisis de perfil latente
<b>Variables observables categóricas</b>	Análisis de rasgos latentes o teoría de respuesta al ítem	Modelos de clases latentes

El objetivo general de realizar un análisis de clases latentes sobre un conjunto de variables observables es llegar a determinar un conjunto de clases latentes que representan los patrones

de respuesta de estas variables y proporcionar una tasa de la prevalencia de cada clase latente (Collins y Lanza; 2009).

## 2.2. Análisis de clases latentes

En esta sección presentamos el análisis de clases latentes de manera más formal y general, tomando como referencia la tesis de Wiener (2015). En ella indica que: tanto las variables observables como latentes son categóricas, y se tienen los siguientes dos supuestos:

1. A nivel de cada clase todos los individuos tienen las mismas probabilidades de respuesta a las variables observables.
2. La independencia condicional; esto es que las respuestas entre los individuos son independientes dado que pertenecen a una misma clase, (página 8).

Wiener (2015) para un mejor entendimiento realiza lo siguiente: supone que hay  $p$  variables observadas dicotómicas; es decir, variables que toman solo dos valores (1 y 0) donde por 1 identificaremos a una respuesta correcta, y sea  $\mathbf{Y}$  la variable latente categórica subyacente con  $c = 1, 2, 3, \dots, C$  clases, que corresponden a las categorías de esta variable. Si denotamos por  $\eta_c = P(Y = c)$  a la probabilidad de que un individuo pertenezca a la clase latente  $c$ , se debe cumplir que:

$$\sum_{c=1}^C \eta_c = 1 .$$

Por otro lado, la probabilidad de responder correctamente a una variable observada  $i$ , condicionada por la pertenencia a la clase latente  $c$ , se denota por  $\pi_{ic}$ . Si  $X_i$  denota a la posible respuesta de un individuo perteneciente a la clase  $c$  a la variable observable  $i$ , esta probabilidad viene dada por:

$$\pi_{ic} = P(X_i = 1 | Y = c) .$$

Asimismo, sea  $\mathbf{X} = (X_1, X_2, \dots, X_p)'$  el vector de respuestas de un individuo a las  $p$  variables observables y sea  $\mathbf{x} = (x_1, x_2, \dots, x_p)'$  el valor observado o patrón de respuesta de este individuo. La probabilidad de observar este patrón particular, condicionada a que el individuo pertenezca a la clase latente  $c$  se obtiene mediante:

$$P(\mathbf{X} = \mathbf{x} | Y = c) = \prod_{i=1}^p \pi_{ic}^{x_i} (1 - \pi_{ic})^{1-x_i} . \quad (2.1)$$

Para determinar la probabilidad de pertenencia de un individuo a una clase latente  $c$ , observado su patrón de respuesta  $\mathbf{x}$ , podemos usar el teorema de Bayes y obtener la probabilidad de clasificación a posteriori como sigue:

$$P(Y = c | \mathbf{X} = \mathbf{x}) = \frac{P(\mathbf{X} = \mathbf{x} | Y = c) \eta_c}{P(\mathbf{X} = \mathbf{x})} . \quad (2.2)$$

donde la función de probabilidad conjunta del vector de respuestas observables o patrones de respuesta  $\mathbf{X}$  viene dado por:

$$P(\mathbf{X} = \mathbf{x}) = \sum_{c=1}^C \eta_c \prod_{i=1}^p \pi_{ic}^{x_i} (1 - \pi_{ic})^{1-x_i} . \quad (2.3)$$

De las ecuaciones (2.1), (2.2) y (2.3) se obtiene la probabilidad a posteriori de que un individuo con un patrón de respuesta  $\mathbf{x} = (x_1, x_2, \dots, x_p)'$  pertenezca a la clase  $c$ , como:

$$P(Y = c | \mathbf{X} = \mathbf{x}) = \frac{\eta_c \prod_{i=1}^p \pi_{ic}^{x_i} (1 - \pi_{ic})^{1-x_i}}{\sum_{c=1}^C \eta_c \prod_{i=1}^p \pi_{ic}^{x_i} (1 - \pi_{ic})^{1-x_i}} . \quad (2.4)$$



Esta función se usa como regla de clasificación para asignar un individuo a la clase con mayor probabilidad de pertenencia, conocido su patrón de respuesta  $\mathbf{x}$ . El principal problema estadístico consiste en la estimación de las probabilidades  $\pi_{ic}$  y las probabilidades de pertenencia a la clases  $\eta_c$ , así como la cuantificación del ajuste de los datos a este modelo. Otro problema será la identificación del número, y la interpretación de las clases latentes subyacentes de manera que tengan un sentido, (página 2).

### 2.3. El modelo de regresión de clases latentes

Según Wiener (2015): el análisis de clases latentes tiene limitaciones, en el sentido que considera que la probabilidad de pertenencia a una clase latente es igual para todos los individuos que tienen los mismos patrones de respuesta a las variables observadas que sirven para obtener la variable subyacente latente, y no toma en cuenta las características propias de cada individuo (sexo, edad, nivel socioeconómico, etc), los cuales podrían influir en su patrón de respuesta, y por ende, en la probabilidad de pertenecer a una clase latente, (página 3).

Una extensión del análisis de clases latentes es el análisis de regresión de clases latentes (LCR), que considera covariables como las señaladas anteriormente, permitiendo de esta manera modelar con mayor información la probabilidad de pertenencia a una clase dependiendo de estas covariables (Clogg; 1981; Formann; 1992).

Este análisis es similar al de la regresión logística, pero con la diferencia de que la variable respuesta es latente en lugar de ser directamente observable. Este evalúa el efecto de un conjunto de covariables en la prevalencia de un individuo a una clase latente, verificando si su inclusión es significativa o no, y en caso de serlo que tipo de efecto genera. Como en cualquier tipo de regresión, el grupo de covariables puede incluir variables categóricas o numéricas (Agresti; 1990).

Sea  $\mathbf{Z} = (Z_1, Z_2, \dots, Z_m)'$  un vector columna  $m$ -dimensional de covariables. Usando el teorema de probabilidad total con la asunción de independencia condicional entre las variables observables, y que se extiende a las covariables se tiene lo siguiente:

$$P(\mathbf{X} = \mathbf{x} | \mathbf{Z} = \mathbf{z}) = \sum_{c=1}^C P(\mathbf{X} = \mathbf{x} | \mathbf{Z} = \mathbf{z}, Y = c) P(Y = c | \mathbf{Z} = \mathbf{z}) ,$$

donde

$$\begin{aligned} P(\mathbf{X} = \mathbf{x} | \mathbf{Z} = \mathbf{z}, Y = c) &= \frac{P(\mathbf{X} = \mathbf{x}, \mathbf{Z} = \mathbf{z}, Y = c)}{P(\mathbf{Z} = \mathbf{z}, Y = c)} = \frac{P(\mathbf{X} = \mathbf{x}, \mathbf{Z} = \mathbf{z} | Y = c) P(Y = c)}{P(\mathbf{Z} = \mathbf{z} | Y = c) P(Y = c)} \\ &= \frac{P(\mathbf{X} = \mathbf{x} | Y = c) P(\mathbf{Z} = \mathbf{z} | Y = c)}{P(\mathbf{Z} = \mathbf{z} | Y = c)} = P(\mathbf{X} = \mathbf{x} | Y = c) \end{aligned} \quad (2.5)$$

Por tanto, la probabilidad de tener un vector de respuestas en las variables observadas dada las covariables queda definido como:

$$P(\mathbf{X} = \mathbf{x} | \mathbf{Z} = \mathbf{z}) = \sum_{c=1}^C P(\mathbf{X} = \mathbf{x} | Y = c) \eta_c(\mathbf{z}) , \quad (2.6)$$

donde  $\eta_c(\mathbf{z}) = P(Y = c | \mathbf{Z} = \mathbf{z})$ .

Cabe señalar que en este modelo, las covariables no están relacionadas con las probabilidades de respuesta a las variables observadas al interior de cada clase latente. Asimismo la ecuación (2.6) es semejante a (2.3), con la diferencia que el parámetro que muestra la probabilidad de pertenencia a la clase se ha definido en función de las covariables. Por tanto, se obtiene la extensión de regresión del modelo de clase latente como sigue:

$$P(\mathbf{X} = \mathbf{x} | \mathbf{Z} = \mathbf{z}) = \sum_{c=1}^C \prod_{i=1}^p \pi_{ic}^{x_i} (1 - \pi_{ic})^{1-x_i} \eta_c(\mathbf{z}) \quad (2.7)$$

en la cual  $\pi_{ic} = P(X_i = 1|Y = c)$  y  $\eta_c(\mathbf{z})$  corresponde a un modelo de regresión logístico multinomial con categoría de referencia 1, expresándose de la siguiente manera:

$$\eta_c(\mathbf{z}) = \frac{e^{\beta_{0c} + \sum_{k=1}^m \beta_{kc} z_k}}{1 + \sum_{l=2}^C e^{\beta_{0l} + \sum_{k=1}^m \beta_{kl} z_k}} = \frac{e^{\tilde{\mathbf{z}}' \boldsymbol{\beta}_c}}{1 + \sum_{l=2}^C e^{\tilde{\mathbf{z}}' \boldsymbol{\beta}_l}}, \quad c = 1, 2, 3, \dots, C \quad (2.8)$$

con  $\tilde{\mathbf{z}}' = [1, \mathbf{z}']$ ,  $\boldsymbol{\beta}_1 = \mathbf{0} \in \mathbb{R}^{m+1}$  y el vector de coeficientes  $\boldsymbol{\beta}_c = (\beta_{0c}, \beta_{1c}, \beta_{2c}, \dots, \beta_{mc})$  perteneciente a la clase latente  $c$ . De esta manera, Wiener (2015): indica que si contamos con  $m$  covariables, el vector  $\boldsymbol{\beta}_c$  será de tamaño  $m+1$ , con un coeficiente por cada covariable más un intercepto. Asimismo, la ecuación (2.8) muestra que el parámetro  $\boldsymbol{\beta}_c$  determina la relación condicional entre covariables y la pertenencia a una clase latente, produciendo un modelo de regresión de clases latentes cuyos parámetros a estimar son los  $C-1$  vectores de coeficientes  $\boldsymbol{\beta}_c$  y las probabilidades  $\pi_{ic}$  de una respuesta correcta condicional de cada variable observable  $i$  a cada clase latente  $c$ , (página 18).

Luego, usando la ecuación (2.5) las distribuciones a posteriori quedan definidas como sigue:

$$\begin{aligned} P(Y = c | \mathbf{X} = \mathbf{x}, \mathbf{Z} = \mathbf{z}) &= \frac{P(Y = c, \mathbf{X} = \mathbf{x}, \mathbf{Z} = \mathbf{z})}{P(\mathbf{X} = \mathbf{x}, \mathbf{Z} = \mathbf{z})} \\ &= \frac{P(\mathbf{X} = \mathbf{x} | Y = c, \mathbf{Z} = \mathbf{z}) P(Y = c, \mathbf{Z} = \mathbf{z})}{P(\mathbf{X} = \mathbf{x}, \mathbf{Z} = \mathbf{z})} \\ &= \frac{P(\mathbf{X} = \mathbf{x} | Y = c) P(Y = c | \mathbf{Z} = \mathbf{z}) P(\mathbf{Z} = \mathbf{z})}{P(\mathbf{X} = \mathbf{x} | \mathbf{Z} = \mathbf{z}) P(\mathbf{Z} = \mathbf{z})} \\ &= \frac{P(\mathbf{X} = \mathbf{x} | Y = c) P(Y = c | \mathbf{Z} = \mathbf{z})}{\sum_{l=1}^C P(\mathbf{X} = \mathbf{x} | Y = l) P(Y = l | \mathbf{Z} = \mathbf{z})} \end{aligned}$$

Finalmente la distribución a posteriori de clasificación para el caso binario viene dado por:

$$P(Y = c | \mathbf{X} = \mathbf{x}, \mathbf{Z} = \mathbf{z}) = \frac{\eta_c(\mathbf{Z}) \prod_{i=1}^p \pi_{ic}^{x_i} (1 - \pi_{ic})^{1-x_i}}{\sum_{l=1}^C \eta_l(\mathbf{Z}) \prod_{i=1}^p \pi_{il}^{x_i} (1 - \pi_{il})^{1-x_i}}. \quad (2.9)$$

En estos modelos de regresión de clases latentes, el desafío es decidir el número de clases a considerar, estimar los parámetros del modelo, e identificar el efecto que produce el incorporar covariables en las probabilidades de pertenencia a una clase.

## 2.4. El modelo regresión de clases latentes con predictores

Otro modelo de regresión de clases latentes y mezcla finita que introduciremos, y abreviaremos por LCRP, es uno de la familia de los modelos lineales generalizados (GLM) en el que los parámetros difieren entre las clases latentes. Una de las aplicaciones de estos modelos es la agrupación en clústers o la segmentación basada en un modelo de regresión. Otra aplicación incluye el modelado de dos niveles. Debido a esta estructura y a las diversas opciones para imponer restricciones a los parámetros, hacen a este modelo el más flexible entre los modelos de clase latente en términos de restricciones.

Se trata de un modelo con una única variable latente nominal  $Y$ , de  $C$  categorías, que se denominan Clústeres o Clases.

Asimismo, el modelo contiene una única variable dependiente  $X_{it}$ , que puede observarse más de una vez para cada caso  $i$ . Estas respuestas múltiples pueden ser replicaciones experimentales, mediciones repetidas en diferentes momentos u ocasiones, observaciones agrupadas, ó respuestas a un conjunto de ítems del cuestionario. El número total de réplicas se denota por  $T_i$ , donde el índice  $i$  en  $T_i$  permite tratar con un número desigual de observaciones por caso.

Asimismo, en este modelo se hace una distinción entre dos tipos de variables exógenas:

- Variables que afectan a la variable latente  $Y$ , que pueden variar entre casos y son usadas para predecir la pertenencia a la clase. Ellas se denominan covariables y son denotados como  $z_{ir}^{cov}$  con  $1 \leq r \leq R$ , siendo  $R$  el número de covariables.
- Variables que afectan a la variable dependiente  $X_{it}$  a través de un GLM, que pueden variar dentro de los casos y se utilizan para predecir las mediciones repetidas de la variable respuesta. Estas se denominan predictores y se denotan como  $z_{itq}^{pred}$  con  $1 \leq q \leq Q$  donde  $Q$  es el número de predictores, y el índice  $t$  en  $z_{itq}^{pred}$  refleja que el valor de un predictor puede cambiar entre réplicas, mientras que una covariable tiene el mismo valor en todas las réplicas de un caso particular.

Por lo anterior, se concluye que el modelo trabaja con un conjunto de datos de dos niveles, donde  $t$  se refiere a las observaciones de nivel inferior que están dentro de las observaciones de nivel superior  $i$ . Las covariables sirven como variables exógenas de nivel superior, y los predictores como variables exógenas de nivel inferior.

Realizando una comparación de este modelo con el LCR, indicado en (2.7) y (2.8), en la cual las covariables no se relacionan con la probabilidades de respuesta a las variables observadas al interior de cada clase latente, se observa que en el nuevo modelo de la familia de los GLM, se plantea un grupo de variables predictoras perteneciente a una clase afectando a la respuesta de un individuo, además dichos predictores pueden cambiar entre réplicas o mediciones repetidas produciendo cambios a su vez en las respuestas. Por otro lado, ambos modelos son similares en lo que respecta a un grupo de covariables afectando a la pertenencia de clase, incluso en el modelo de la familia de los GLM estas covariables no cambiarán en los casos en que haya réplicas para una respuesta en particular.

#### 2.4.1. Estructura de la probabilidad y predictores lineales

La estructura de la probabilidad del modelo LCRP tiene la siguiente forma para la función de probabilidad y/o densidad de la  $i$ -ésima variable dependiente:

$$f(\mathbf{x}_i | \mathbf{z}_i^{cov}, \mathbf{z}_i^{pred}) = \sum_{c=1}^C P(Y = c | \mathbf{z}_i^{cov}) \prod_{t=1}^{T_i} f(x_{it} | Y = c, \mathbf{z}_{it}^{pred}) . \quad (2.10)$$

En la cual se hace la distinción entre covariables y predictores, además se permite diferentes números de réplicas para cada caso y se asume que la función de masa condicional  $f(x_{it} | Y = c, \mathbf{z}_{it}^{pred})$  tienen la misma forma para cada  $t$  y no se permite efectos entre las múltiples respuestas.

En el caso particular que  $X_{it}$  es una variable nominal que toma el valor de 1 si se da una respuesta correcta y 0 en caso contrario, entonces la función de distribución para cada  $X_{it}$  es de la forma:

$$P(X_{it} = 1 | Y = c, \mathbf{z}_{it}) = \pi_{ict} = \frac{\exp(\eta_{ict})}{1 + \exp(\eta_{ict})} , \quad (2.11)$$

donde  $\pi_{ict}$  denota a la probabilidad de dar una respuesta correcta, dado que el individuo  $i$  pertenece a la clase  $c$  y tiene en  $t$  el conjunto de predictores  $\mathbf{z}_{it}$ , asimismo  $\eta_{ict}$  es un predictor lineal de la forma:

$$\eta_{ict} = \beta_{c10} + \sum_{q=1}^Q \beta_{c1q} \cdot z_{itq}^{pred} , \quad (2.12)$$

donde  $\beta_{c10}$  es un parámetro de intercepto cuando la respuesta toma el valor de 1, y los  $\beta_{c1q}$  son coeficientes de regresión para esa misma respuesta correspondiente al  $q$ -ésimo predictor.



## 2.5. El modelo DINA

Los modelos de diagnóstico cognitivo (CDM) se utilizan para estimar las habilidades latentes que los examinados pudiesen tener para responder correctamente a los ítems de una prueba. Uno de los modelos más populares es el modelo DINA.

Según DeCarlo (2011) la idea básica del modelo DINA es que para responder correctamente un ítem, el examinado debe de poseer todas las habilidades diseñadas para este ítem. Este modelo plantea un contexto en el que  $N$  examinados denotados por  $i = 1, 2, \dots, N$ , son evaluados mediante una prueba con  $J$  ítems denotados por  $j = 1, 2, \dots, J$ , cada uno de los cuales requiere la presencia de  $K_j$  atributos para poder ser correctamente respondidos, donde  $K_j \leq K$ , siendo  $K$  la totalidad de habilidades evaluados por la prueba.

El modelo DINA requiere la construcción de una matriz  $Q = [q_{jk}]$  (Tatsuoka; 1985). Esta matriz de orden  $J \times K$  indica las habilidades requeridas para responder correctamente cada ítem  $j$ , donde  $k = 1, 2, \dots, K$ :

$$Q = \begin{pmatrix} q_{11} & q_{12} & \cdots & q_{1k} \\ q_{21} & q_{22} & \cdots & q_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ q_{J1} & q_{J2} & \cdots & q_{JK} \end{pmatrix}.$$

Ella es una matriz binaria que especifica la relación entre los ítems y las habilidades; en la cual un valor de 1 indica que el atributo particular es requerido para responder el ítem, mientras que un valor de 0 representa que la habilidad particular no es necesaria. Así, cada elemento  $q_{jk}$  de la matriz  $Q$  se define como:

$$q_{jk} = \begin{cases} 1 & , \text{cuando la habilidad } k \text{ se requiere para responder correctamente el ítem } j \\ 0 & , \text{en caso contrario} \end{cases}$$

De esta manera, cada fila de la matriz  $Q$  consta de ceros y unos, en la cual un valor de cero indica que no se necesita de la habilidad  $k$  para responder el ítem  $j$ , y un valor de uno que indica que se requiere de esa habilidad.

Tomando como referencia la tesis de Sosa (2017), en ella indica que: el modelo DINA estima la relación entre los  $N$  individuos y las  $K$  habilidades basándose en el patrón de respuestas de los individuos a los ítems, que se expresan en una matriz de unos y ceros, donde se asigna 1 cuando el examinado  $i$  responde correctamente a un ítem  $j$  y 0 en caso contrario, para lo cual se define una matriz de orden  $N \times K$ ,  $A$ , como:

$$A = \begin{pmatrix} \alpha_{11} & \alpha_{12} & \cdots & \alpha_{1k} \\ \alpha_{21} & \alpha_{22} & \cdots & \alpha_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ \alpha_{N1} & \alpha_{N2} & \cdots & \alpha_{NK} \end{pmatrix},$$

la cual posee, igual que  $Q$ , una estructura binaria; es decir, con elementos 0 y 1, donde cada  $\alpha_{ik}$  se define como:

$$\alpha_{ik} = \begin{cases} 1 & , \text{si el examinado } i \text{ domina la habilidad } k \\ 0 & , \text{en caso contrario} \end{cases}$$

Sin embargo, a diferencia de  $Q$ , las entradas de esta matriz son latentes; es decir, no observables. Además se define el vector latente  $\alpha_i = (\alpha_{i1}, \alpha_{i2}, \dots, \alpha_{iK})$  que indica todas las

habilidades que domina o no el examinado  $i$ . Estos que denominaremos perfiles latentes, son desconocidos y su estimación es uno de los objetivos del modelo DINA (página 7). Este vector es llamado también el “estado de conocimiento” del examinado  $i$  (Junker y Sijtsma; 2001).

El modelo DINA es un modelo conjuntivo debido a que para que un examinado responda correctamente un ítem, debe tener todas las habilidades requeridas por el ítem (Maris; 1999). Sosa (2017) indica que: mediante la variable latente binaria  $\eta_{ij}$ , que clasifica si el examinado  $i$  domina todos los atributos requeridos para el ítem  $j$  (página 8), y se determina como:

$$\eta_{ij} = \prod_{k=1}^K \alpha_{ik}^{q_{jk}} , \quad (2.13)$$

donde

$$\eta_{ij} = \begin{cases} 1 & , \text{ si el examinado } i \text{ domina todas las habilidades que se necesitan} \\ & \text{ para responder correctamente el ítem } j, \\ 0 & , \text{ en caso contrario} \end{cases}$$

Por tanto, este sería la entrada determinística del modelo, en la cual se calcula la presencia o ausencia del conjunto necesario de habilidades haciendo que  $\eta_{ij}$  sea cero o uno. Como se ha indicado, se trata de un modelo conjuntivo en la cual un examinado  $i$  debe tener todas las habilidades necesarias para resolver el problema. Supongamos, por ejemplo, que de un conjunto de cuatro habilidades, la 1 y 3 son necesarias para resolver el primer ítem, entonces la ecuación (2.13) junto con la matriz  $Q$  nos da como resultado para el primer ítem,  $\eta_{i1} = \alpha_{i1}^1 \times \alpha_{i2}^0 \times \alpha_{i3}^1 \times \alpha_{i4}^0 = \alpha_{i1} \times \alpha_{i3}$ . Por lo tanto, debido al término multiplicativo,  $\eta_{ij}$  es cero si  $\alpha_{i1}$  o  $\alpha_{i3}$ , o ambos, son cero, y  $\eta_{ij}$  es igual a 1 si y solo si ambas habilidades están presentes, esto vendría a ser el aspecto conjuntivo del modelo.

Sea  $Y_{ij}$  una variable aleatoria correspondiente a la respuesta del examinado  $i$ , al ítem  $j$ . El modelo DINA calcula la probabilidad de que el examinado  $i$  resuelva el ítem  $j$  correctamente dada su capacidad como,

$$p(Y_{ij} = 1 | \eta_{ij}) = (1 - s_j)^{\eta_{ij}} (g_j)^{1 - \eta_{ij}} . \quad (2.14)$$

En (2.14) se introducen dos parámetros: el desliz  $s_j$  y la adivinación  $g_j$  (Junker y Sijtsma; 2001). Los estudiantes que poseen todos los atributos requeridos para un ítem, pueden tener un desliz y responder incorrectamente al ítem, mientras que los estudiantes que no poseen todos los atributos requeridos para un ítem, pueden adivinar y responderlo correctamente. El modelo DINA define el parámetro de desliz como  $s_j = p(Y_{ij} = 0 | \eta_{ij} = 1)$  y el parámetro de adivinación como  $g_j = p(Y_{ij} = 1 | \eta_{ij} = 0)$  para cada ítem  $j$ .

Si usamos la terminología de la teoría de detección de señales (Macmillan y Creelman; 2005),  $g_j$  viene a ser la probabilidad de una falsa alarma,  $1 - s_j$  la probabilidad de un golpe y  $s_j$  es la probabilidad de una falla. Se debe tener en cuenta que los parámetros de adivinación y desliz están condicionados a las habilidades específicas indicadas en la matriz  $Q$ , por lo que cualquier cambio en ella puede significar cambiar los parámetros de desliz y adivinación.

## 2.6. El modelo RDINA

El modelo DINA se puede ajustar dentro de un marco bayesiano utilizando los métodos de cadenas de Markov Monte Carlo (MCMC) (De La Torre y Douglas; 2004, 2008; Henson et al.; 2009; Junker y Sijtsma; 2001). Sin embargo, DeCarlo (2011) propone también ajustar el modelo utilizando la estimación de máxima verosimilitud (MLE) para lo cual realiza una reparametrización del modelo DINA como un modelo de regresión logística con clases latentes.

Esta reparametrización consiste en escribir el parámetro de adivinación  $g_j$  en términos de un nuevo parámetro  $f_j$  mediante,

$$g_j = p(Y_{ij} = 1 | \eta_{ij} = 0) = \exp(f_j) / (1 + \exp(f_j)) \quad (2.15)$$

Lo anterior, que nos brinda la probabilidad de que un examinado responda correctamente un ítem dado que no tiene las habilidades requeridas (falsa alarma), se puede escribir en forma más simple usando una función *logit*,  $\text{logit}(p) = \log(p/(1 - p))$ , como :

$$\text{logit} [p(Y_{ij} = 1 | \eta_{ij} = 0)] = f_j \quad (2.16)$$

El parámetro  $f_j$  representa entonces el log odds de una falsa alarma (probabilidad que el sujeto adivine la respuesta). Similarmente, [DeCarlo \(2011\)](#) reescribe el *log odds* de que un examinado responda correctamente a un ítem  $i$  dado que tiene las habilidades requeridas (hit) como:

$$\text{logit} [p(Y_{ij} = 1 | \eta_{ij} = 1)] = f_j + d_j \quad (2.17)$$

Las dos ecuaciones anteriores corresponden a una reparametrización del modelo DINA (2.13), y ambas se pueden escribir en un solo modelo de la siguiente manera:

$$\text{logit} [p(Y_{ij} = 1 | \eta_{ij})] = f_j + d_j \eta_{ij} \quad (2.18)$$

La ecuación (2.18) representa entonces un modelo logístico de clases latentes, en la cual los ítems sirven como detectores de los conjuntos de habilidades. Dado que esta ecuación se trata simplemente de una reparametrización del modelo DINA, esta se denomina un modelo reparametrizado DINA (RDINA, reparameterized deterministic input noisy and gate), [DeCarlo \(2011\)](#). Aquí como dijimos, el parámetro  $f_j$  proporciona las estimaciones de *log odds* de una falsa alarma, mientras que el parámetro  $d_j$  proporciona una medida de qué tan bien el ítem detecta la presencia frente a la ausencia del conjunto de habilidades requerido, es decir,  $d_j$  ayuda a discriminar entre los examinados con y sin las habilidades requeridas.

Exponenciando estos parámetros, se puede recuperar el parámetro de desliz mediante:

$$\begin{aligned} \text{logit} p(Y_{ij} = 1 | \eta_{ij} = 1) &= f_j + d_j \\ \text{logit} (1 - p(Y_{ij} = 0 | \eta_{ij} = 1)) &= f_j + d_j \\ \text{logit} (1 - s_j) &= f_j + d_j \\ \log\left(\frac{1 - s_j}{s_j}\right) &= f_j + d_j \\ \left(\frac{1 - s_j}{s_j}\right) &= \exp(f_j + d_j) \\ s_j &= 1 / (1 + \exp(f_j + d_j)) \\ s_j &= 1 - \exp(f_j + d_j) / (1 + \exp(f_j + d_j)) \end{aligned} \quad (2.19)$$

Con los parámetros  $f_j$  y  $d_j$ , es posible también calcular en el RDINA un índice de discriminación ([De La Torre; 2008](#)), que representa qué tan bien un ítem es capaz de clasificar a un examinado por haber dominado una habilidad para cada ítem. Este índice tiene la siguiente forma:

$$\gamma_j = 1 - g_j - s_j = \exp(f_j + d_j) / (1 + \exp(f_j + d_j)) - \exp(f_j) / (1 + \exp(f_j)) \quad (2.20)$$

Por ejemplo, una disminución en la media de las estimaciones del índice de discriminación en el modelo RDINA con covariable afectando a las habilidades al pasar a un modelo RDINA con covariable afectando a los ítems, podría ser debido a las estimaciones más bajas en el parámetro de adivinación y más altos en el de desliz; la medida en que los parámetros de desliz cambiaron fue mayor que la medida en que los parámetros de adivinación disminuyeron.

## Capítulo 3

# Extensión al modelo DINA con covariable

Los CDM se desarrollaron para proporcionar información más específica en forma de perfiles que resuelven la limitación de los métodos clásicos y los modelos de teoría de respuesta al ítem unidimensionales (IRT). Varios CDMs han sido propuestos en la literatura, proporcionando un marco ideal para realizar un análisis para clasificar a los examinados en perfiles de habilidades que indican su dominio en ella. En este estudio ha sido seleccionado el modelo DINA dada su parsimonia, facilidad de interpretación y posible extensión a modelos de diagnósticos cognitivos más complejos. En nuestro caso específico, ampliaremos el marco del modelo DINA al incluir una covariable en base a su modelo reparameterizado RDINA.

Park y Lee (2014) investigaron cómo la extensión con covariable en el modelo RDINA se puede aplicar a datos del mundo real en alumnos del cuarto grado, y realizaron un estudio con la finalidad de examinar el efecto de sus capacidades en ciencia, en este caso del puntaje correcto en la evaluación de la ciencia TIMSS (Tendencias en el Estudio Internacional de Matemáticas y Ciencias) podría tener sobre la probabilidad de un examinado para resolver ítems de matemáticas; encontrando en los resultados que la capacidad científica tuvo un efecto significativo tanto en los ítems como en las habilidades. Los estudiantes con mayor capacidad científica, mostraron tener una mayor probabilidad de resolver correctamente los ítems matemáticos y de ser clasificados como dominantes en seis de los siete habilidades especificados en la matriz  $Q$ . Asimismo, los resultados de simulación que realizaron mostraron una recuperación estable de parámetros y clases latentes para diferentes tamaños de muestra. Estos hallazgos no sugieren una relación causal entre los puntajes de ciencias y la capacidad matemática; pero ayudan a comprender cómo diversos factores influyen en el dominio de las habilidades. Además, la asociación significativa que se encontró para estas habilidades puede conducir a estudios adicionales con resultados significativos para los investigadores aplicados.

### 3.1. El modelo RDINA con covariable

En el modelo RDINA o DINA reparametrizado, las covariables se pueden especificar en dos niveles: en el nivel inferior, sobre el cual afectan a la forma en que los examinados resuelven los ítems (es decir, la probabilidad de respuesta), y en el nivel superior, en la que influyen en el dominio de los atributos (es decir, la clasificación latente).

Al tratar el modelo DINA como un modelo de clase latente, es necesario considerar que los datos básicos para cada examinado  $i$  son sus patrones de respuesta observados  $\mathbf{y}_i = (y_{i1}, y_{i2}, \dots, y_{iJ})$  para los  $J$  ítems de una prueba (es decir, si responden correctamente o no a los ítems de la prueba), siendo sus probabilidades de presentar dicho patrón  $P(Y_{i1} = y_{i1}, Y_{i2} = y_{i2}, \dots, Y_{iJ} = y_{iJ})$ , que también puede denotarse más compactamente como  $P(\mathbf{Y}_i = \mathbf{y}_i)$ . Sea, por otro lado,  $\boldsymbol{\alpha}_c = (\alpha_{c1}, \alpha_{c2}, \dots, \alpha_{cK})$  el estado de conocimiento de cualquier examinado en

la clase  $c$  para las  $K$  habilidades que busca medir la prueba, entendiéndose que todos los examinados en ella poseen idéntico estado de conocimiento y que  $\alpha_{cK}$  toma el valor de 1 si estos examinados poseen la habilidad  $k$ . Aquí, se considera a  $C$  como el número de clases latentes con valor máximo posible  $C = 2^K$  si todo posible patrón es admisible.

El modelo de clase latente, relaciona las probabilidades del patrón de respuesta con las probabilidades de respuesta condicional de la ecuación (2.18) de la siguiente manera (DeCarlo; 2011):

$$P(Y_{i1} = y_{i1}, Y_{i2} = y_{i2}, \dots, Y_{iJ} = y_{iJ}) = \sum_{c=1}^C p(\alpha_c) P(Y_{i1} = y_{i1}, Y_{i2} = y_{i2}, \dots, Y_{iJ} = y_{iJ} | \alpha_c), \quad (3.1)$$

donde  $p(\alpha_c)$  es la probabilidad de que el examinado  $i$  pertenezca a la clase  $c$ .

Lo anterior simplemente muestra una premisa básica del análisis de clases latentes: que la probabilidad no condicional del vector de respuesta es una suma ponderada de las probabilidades condicionales sobre las clases latentes. Aplicando el supuesto básico de independencia condicional para las respuestas a las clases latentes (Clogg; 1995), se tiene que:

$$P(Y_{i1} = y_{i1}, Y_{i2} = y_{i2}, \dots, Y_{iJ} = y_{iJ} | \alpha_c) = \prod_{j=1}^J P(Y_{ij} = y_{ij} | \alpha_c). \quad (3.2)$$

Por lo tanto, de (3.1) y (3.2) la probabilidad de respuesta del examinado  $i$  a los  $J$  ítems se puede modelar como:

$$P(Y_{i1} = y_{i1}, Y_{i2} = y_{i2}, \dots, Y_{iJ} = y_{iJ}) = \sum_{c=1}^C p(\alpha_c) \prod_{j=1}^J P(Y_{ij} = y_{ij} | \alpha_c), \quad (3.3)$$

donde el término  $P(Y_{ij} = y_{ij} | \alpha_c)$  se obtiene del modelo DINA y la ecuación (2.14), o del modelo RDINA y la ecuación (2.18). Utilizándose en esta última la función de enlace *logit* para la probabilidad de respuesta al ítem  $j$  de un examinado  $i$  en la clase  $c$ , se tiene que:

$$\text{logit } P(Y_{ij} = y_{ij} | \alpha_c) = (f_j + d_j \eta_{cj})^{y_{ij}} (1 - (f_j + d_j \eta_{cj}))^{1-y_{ij}}, \quad (3.4)$$

donde  $\eta_{cj} = \prod_{k=1}^K \alpha_{ck}^{q_{jk}}$ .

Se requiere, por otro lado, asumir una estructura para  $p(\alpha_c)$  para lo cual aplicando el supuesto restrictivo de independencia en la estructura de habilidades,  $p(\alpha_c)$  se convierte en:

$$p(\alpha_c) = p(\alpha_{c1}, \alpha_{c2}, \dots, \alpha_{cK}) = \prod_{k=1}^K p(\alpha_{ck}), \quad (3.5)$$

donde  $p(\alpha_{ck})$  es la probabilidad de que un examinado en la clase  $c$  posea la habilidad  $k$ .

Usando como modelo el logístico para  $p(\alpha_{ck})$  se obtiene luego que:

$$p(\alpha_{ck}) = \exp(b_k) / (1 + \exp(b_k)), \quad (3.6)$$

donde  $b_k$  es un parámetro de “facilidad” para la habilidad  $k$ .

Park y Lee (2014) introducen en este contexto una covariable discreta o continua,  $Z$  para la probabilidad de respuesta (3.3) y que es propia del examinado  $i$ . Usando el teorema de probabilidad total al condicionar a la clase de pertenencia del examinado, se propone entonces el siguiente modelo para la probabilidad de que un examinado  $i$  presente un patrón de respuestas  $\mathbf{y} = (y_1, y_2, \dots, y_J)$ :

$$P(Y_{i1} = y_1, Y_{i2} = y_2, \dots, Y_{iJ} = y_J | Z_i = z) = P(\mathbf{Y}_i = \mathbf{y} | Z_i = z) =$$



$$\sum_{c=1}^C p(\alpha_c | Z_i = z) P(\mathbf{Y}_i = \mathbf{y} | \alpha_c, Z_i = z) ,$$

donde  $P(\mathbf{Y}_i = \mathbf{y} | \alpha_c, Z_i = z)$  es una probabilidad de respuesta similar a (3.4) pero que ahora podría depender del valor de la covariable  $Z_i$ . Aplicándose el supuesto básico de independencia condicional se tiene que:

$$P(\mathbf{Y}_i = \mathbf{y} | \alpha_c, Z_i = z) = \prod_{j=1}^J P(Y_{ij} = y_j | \alpha_c, Z_i = z) . \quad (3.7)$$

Por tanto la probabilidad de respuesta del examinado  $i$  a los  $j$  ítems condicionado al valor de su covariable viene dada por:

$$P(Y_{i1} = y_1, Y_{i2} = y_2, \dots, Y_{iJ} = y_J | Z_i = z) = \sum_{c=1}^C p(\alpha_c | Z_i = z) \prod_{j=1}^J P(Y_{ij} = y_j | \alpha_c, Z_i = z) . \quad (3.8)$$

Nótese que la covariable afecta a las probabilidades de respuesta,  $P(Y_{ij} = y_{ij} | \alpha_c, Z_i = z)$ . La siguiente ecuación modela estas afectaciones como:

$$\text{logit } P(Y_{ij} = y_{ij} | \alpha_c, Z_i = z) = (f_j + d_j \eta_{cj} + \ell_j z)^{y_{ij}} (1 - (f_j + d_j \eta_{cj} + \ell_j z))^{1-y_{ij}} , \quad (3.9)$$

donde  $\ell_j$  es un parámetro de regresión del ítem  $j$  sobre la covariable  $Z_i$ . De otro lado, la covariable se asume que también afecta a las probabilidades de pertenencia a las clases,

$$p(\alpha_c | Z_i = z) = \prod_{k=1}^K p(\alpha_{ck} | Z_i = z) , \quad (3.10)$$

mediante el modelo de regresión logística,

$$p(\alpha_{ck} | Z_i = z) = \exp(b_k + h_k z) / (1 + \exp(b_k + h_k z)) , \quad (3.11)$$

siendo  $h_k$  un parámetro de regresión de la habilidad  $k$  sobre la covariable  $Z_i$ .

Cuando la covariable condiona las probabilidades de respuesta, esta puede cambiar la estimación de los parámetros de adivinación y desliz. Por ejemplo, si asumimos que  $Z_i$  es una covariable binaria, entonces la presencia o la ausencia de ella (es decir,  $z = 1$  ó  $z=0$ ) puede influir en las tasas de adivinación y desliz en el modelo DINA por un factor que depende de  $\ell_j$ , condicional al valor de la covariable. De esta manera, la ecuación (3.12) muestra el cambio en el parámetro de adivinación de los examinados cuando pasan de  $z = 0$  ( $g_{j0}$ ) a  $z = 1$  ( $g_{j1}$ ), y la ecuación (3.13) muestra el cambio en el parámetro de desliz de los examinados cuando pasan de  $z = 0$  ( $s_{j0}$ ) a  $z = 1$  ( $s_{j1}$ ). Para covariables continuas, el parámetro  $\ell_j$  indica cambios en estos dos parámetros para un incremento unitario en  $z$ :

$$\begin{aligned} g_{j1} - g_{j0} &= \frac{\exp(f_j + \ell_j)}{1 + \exp(f_j + \ell_j)} - \frac{\exp(f_j)}{1 + \exp(f_j)} \\ &= \exp(f_j) \left( \frac{\exp(\ell_j) - 1}{(1 + \exp(f_j + \ell_j))(1 + \exp(f_j))} \right) . \end{aligned} \quad (3.12)$$

$$\begin{aligned} s_{j1} - s_{j0} &= \left( 1 - \frac{\exp(f_j + d_j + \ell_j)}{1 + \exp(f_j + d_j + \ell_j)} \right) - \left( 1 - \frac{\exp(f_j + d_j)}{1 + \exp(f_j + d_j)} \right) \\ &= \exp(f_j + d_j) \left( \frac{1 - \exp(\ell_j)}{(1 + \exp(f_j + d_j + \ell_j))(1 + \exp(f_j + d_j))} \right) . \end{aligned} \quad (3.13)$$

Cuando la covariable condiciona las probabilidades de las habilidades indicada en la ecuación (3.11), ésta se convierte en un predictor de los patrones de las habilidades que afectan a la membresía de la clase latente. Similar a la interpretación del parámetro a nivel de ítem,  $\ell_j$ , el parámetro  $h_k$  se relaciona al cambio en el parámetro de facilidad de la habilidad,  $b_k$ , cuando la covariable está presente. Siguiendo el ejemplo con  $Z$  como covariable binaria, la ecuación (3.14) muestra el cambio que se produce en la probabilidad de que un examinado de la clase  $c$  posee la habilidad  $k$  cuando pasa de  $z = 0$  ( $p_0(\alpha_{ck})$ ) a  $z = 1$  ( $p_1(\alpha_{ck})$ ):

$$\begin{aligned} p_1(\alpha_{ck}) - p_0(\alpha_{ck}) &= \frac{\exp(b_k + h_k)}{1 + \exp(b_k + h_k)} - \frac{\exp(b_k)}{1 + \exp(b_k)} \\ &= \exp(b_k) \left( \frac{\exp(h_k) - 1}{(1 + \exp(b_k + h_k))(1 + \exp(b_k))} \right). \end{aligned} \quad (3.14)$$

Comparando el modelo RDINA con covariable con el modelo de regresión de clases latentes se tiene las siguientes diferencias:

- La covariable en el modelo de regresión de clases latentes, afecta a la probabilidad de pertenencia a la clase de la misma manera que en el RDINA, pero en el caso de éste afecta específicamente a la probabilidad de las habilidades, convirtiéndose de esta manera en un predictor de los patrones de las habilidades que afectan a su vez a la membresía de la clase latente.
- La covariable en algunos modelos de regresión de clases latentes, no se relaciona con las probabilidades de respuesta a las variables observadas al interior de cada clase latente, en cambio en el modelo RDINA con covariable, esta última sí condiciona las probabilidades de respuesta al interior, cambiando la estimación de sus parámetros de adivinación y desliz.
- En el modelo de regresión de clases latentes, se requiere la asunción de independencia condicional entre las variables observadas que se extienden a las covariables, mientras que en el modelo RDINA con covariable, se requiere la asunción de independencia condicional para las respuestas a las clases latentes y las covariables.

## Capítulo 4

# Estimación del modelo

El modelo RDINA con covariable puede verse como una analogía de la extensión del análisis de regresión de clases latentes con predictores (LCRP). Para la estimación de este modelo desarrollaremos el método de máxima verosimilitud y el método de la moda a posteriori, vía el algoritmo de Esperanza-Maximización (EM) y de Newton-Raphson. La implementación de este se encuentra en el paquete Latent Gold, por [Vermunt y Magidson \(2015, 2016\)](#).

### 4.1. Función de verosimilitud

Utilizando el modelo (3.8) y escribiendo al patrón de respuestas del individuo  $i$  como  $\mathbf{y}_i = (y_{i1}, y_{i2}, \dots, y_{iJ})$  y el patrón de covariables de los  $N$  individuos como  $\mathbf{z} = (z_1, z_2, \dots, z_N)$ , se tiene que la función de verosimilitud al observar una muestra de las respuestas de  $N$  individuos a los  $J$  ítems está dada por:

$$\mathbf{L}(\boldsymbol{\theta}; \mathbf{y}, \mathbf{z}) = \prod_{i=1}^N P(\mathbf{Y}_i = \mathbf{y}_i | Z_i = z_i) = \prod_{i=1}^N \sum_{c=1}^C p(\boldsymbol{\alpha}_c | Z_i = z_i) \prod_{j=1}^J P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_c, Z_i = z_i), \quad (4.1)$$

donde, aplicándose el supuesto de independencia condicional, la probabilidad de pertenencia del individuo  $i$  a la clase  $c$  condicionado al valor de su covariable  $Z_i$  viene dada por

$$p(\boldsymbol{\alpha}_c | Z_i = z_i) = \prod_{k=1}^K p(\alpha_{ck} | Z_i = z_i),$$

con el término  $p(\alpha_{ck} | Z_i = z)$  dado en (3.11).

Asimismo, la probabilidad de respuesta del examinado  $i$  al ítem  $j$  condicionado a la membresía a la clase  $c$  y a la covariable viene dada por

$$P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_c, Z_i = z_i) = \left( \frac{e^{f_j + d_j \eta_{cj} + \ell_j z_i}}{1 + e^{f_j + d_j \eta_{cj} + \ell_j z_i}} \right)^{y_{ij}} \left( \frac{1}{1 + e^{f_j + d_j \eta_{cj} + \ell_j z_i}} \right)^{1 - y_{ij}}.$$

En (4.1)  $\mathbf{y}$  denota a un arreglo de todas las respuestas de cada uno de los  $N$  individuos a la prueba.

Introduciéndose las funciones de enlace respectivas dentro del modelo, la función de verosimilitud (4.1) queda explícitamente caracterizada por:

$$\mathbf{L}(\boldsymbol{\theta}; \mathbf{y}, \mathbf{z}) = \prod_{i=1}^N P(\mathbf{Y}_i = \mathbf{y}_i | Z_i = z_i) =$$



$$\prod_{i=1}^N \sum_{c=1}^C \prod_{k=1}^K \left( \frac{e^{b_k + h_k z_i}}{1 + e^{b_k + h_k z_i}} \right) \prod_{j=1}^J \left( \left( \frac{e^{f_j + d_j \eta_{cj} + \ell_j z_i}}{1 + e^{f_j + d_j \eta_{cj} + \ell_j z_i}} \right)^{y_{ij}} \left( \frac{1}{1 + e^{f_j + d_j \eta_{cj} + \ell_j z_i}} \right)^{1-y_{ij}} \right), \quad (4.2)$$

donde  $\mathbf{Y}_i = (Y_{i1}, Y_{i2}, \dots, Y_{iJ})$ ,  $\mathbf{y}_i = (y_{i1}, y_{i2}, \dots, y_{iJ})$  y el vector de parámetros a estimar es:

$$\boldsymbol{\theta} = (\mathbf{f}, \mathbf{d}, \boldsymbol{\ell}, \mathbf{b}, \mathbf{h})$$

siendo

- $\mathbf{f} = (f_1, f_2, \dots, f_J)$  un vector de parámetros de los *log odds* de una *falsa alarma*, es decir, de las probabilidades de respuesta correcta a cada ítem por parte de examinados que no tienen las habilidades requeridas.
- $\mathbf{d} = (d_1, d_2, \dots, d_J)$  un vector de parámetros cuyos elementos representan una medida de qué tan bien cada ítem discrimina entre los examinados que dominan todas las habilidades necesarias para responder correctamente este ítem de aquellos que solo dominan algunas de tales habilidades requeridas o ninguna de ellas.
- $\boldsymbol{\ell} = (\ell_1, \ell_2, \dots, \ell_J)$  un vector de parámetros de regresión cuyos elementos representan los cambios en los parámetros de adivinación y desliz de cada ítem para un incremento unitario de la covariable.
- $\mathbf{b} = (b_1, b_2, \dots, b_K)$  un vector de parámetros de regresión cuyos elementos representan una medida de “facilidad” para cada uno de los atributos, dado que determinan la probabilidad de que un individuo posea cada una de las habilidades.
- $\mathbf{h} = (h_1, h_2, \dots, h_K)$  el vector de parámetros cuyos elementos representan los cambios en la probabilidad de que un individuo posea cada una de las habilidades para un incremento unitario de la covariable.

Por tanto la función de log verosimilitud queda representada por:

$$\mathcal{L}(\boldsymbol{\theta}; \mathbf{y}, \mathbf{z}) = \log(\mathbf{L}(\boldsymbol{\theta}; \mathbf{y}, \mathbf{z})) = \sum_{i=1}^N \log \left( \sum_{c=1}^C \prod_{k=1}^K \left( \frac{e^{b_k + h_k z_i}}{1 + e^{b_k + h_k z_i}} \right) \prod_{j=1}^J \left( \left( \frac{e^{f_j + d_j \eta_{cj} + \ell_j z_i}}{1 + e^{f_j + d_j \eta_{cj} + \ell_j z_i}} \right)^{y_{ij}} \left( \frac{1}{1 + e^{f_j + d_j \eta_{cj} + \ell_j z_i}} \right)^{1-y_{ij}} \right) \right) \quad (4.3)$$

## 4.2. Identificabilidad

En este tipo de modelos se puede presentar un posible problema que es la no identificabilidad local de los parámetros, en la que podemos señalar que el modelo es identificable localmente para un valor en específico  $\boldsymbol{\theta}^*$  del vector de parámetros, cuando la función de log-verosimilitud es determinada solamente por los parámetros en una vecindad de  $\boldsymbol{\theta}^*$  (Vermunt y Magidson; 2016; Chung; 2003).

Una forma de evaluar la identificabilidad local del modelo es evaluar si la matriz de información de Fisher posee valores propios mayores que 0. Formann (1985, 1992) mostró que este enfoque es equivalente a evaluar el rango de la matriz jacobiana.

Una práctica estándar para verificar la identificabilidad es utilizar múltiples conjuntos de valores iniciales para la estimación de parámetros. Los diferentes conjuntos de valores iniciales que producen el mismo máximo de probabilidad deberían dar como resultado las mismas estimaciones de los parámetros finales, si esto no ocurre, el modelo no es identificable.

En estos modelos pueden surgir complicaciones en la aplicación de la matriz de información de Fisher para un análisis dado. Lo ideal sería determinar aquellas regiones del espacio de parámetros en las que el modelo dado es localmente identificable. Debido a que esto suele ser difícil desde el punto de vista computacional, estos métodos a menudo se evalúan con respecto a los parámetros estimados para establecer la identificabilidad local del modelo en los valores estimados (Goodman; 1974).

### 4.3. Estimación

#### 4.3.1. Estimación de parámetros por el método de la moda a posteriori

La estimación de parámetros puede realizarse mediante el método de la maximización de la moda a posteriori o también llamado de maximización a posteriori (MAP). Este es el método que se encuentra implementado en el software Latent Gold y el cual utilizaremos en este trabajo.

El método MAP involucra el uso de una distribución a priori para  $\theta$ , que denotaremos por  $p(\theta)$ , y de la distribución a posteriori  $\mathbf{G}(\theta)$ . La estimación MAP implica encontrar los valores de  $\theta$  que maximicen la distribución a posteriori; es decir, que maximicen

$$\begin{aligned}\log \mathbf{G}(\theta) &= \log \mathbf{L}(\theta; \mathbf{y}, \mathbf{z}) + \log p(\theta) \\ &= \sum_{i=1}^N \log P(\mathbf{Y}_i = \mathbf{y}_i | Z_i = z) + \log p(\theta),\end{aligned}\tag{4.4}$$

la cual se podría obtener encontrando el punto donde  $\frac{\partial \log \mathbf{G}(\theta)}{\partial \theta} = 0$ . Nótese que los hiperparámetros definidos en la distribución a priori  $p(\theta)$  podrían elegirse de tal manera que  $\log p(\theta) = 0$ , haciendo que la estimación MAP se convierta en una estimación de máxima verosimilitud (MV). La estimación MAP puede verse también como una estimación de MV penalizada, en la que  $p(\theta)$  sirve como una función que penaliza soluciones que están demasiado cerca de la frontera del espacio paramétrico y, por lo tanto, suaviza las estimaciones lejos de la frontera. De (4.4) se sigue que el método de MV es un caso especial de la maximización a posteriori (MAP), cuando la priori sigue una distribución uniforme o impropia. Esta sería la conexión entre los métodos MAP y MV.

#### 4.3.2. El algoritmo EM

Para encontrar las estimaciones MAP o MV de  $\theta$  en el modelo, el software Latent Gold hace uso del algoritmo EM y de Newton-Raphson. El proceso de estimación comienza con una serie de iteraciones EM y cuando está lo suficientemente cerca de la solución final el programa cambia al método de Newton-Raphson. De esta manera, se explota las ventajas de ambos algoritmos; es decir, la estabilidad de EM incluso cuando está lejos del óptimo y la velocidad de Newton-Raphson cuando está cerca del óptimo.

Para obtener las estimaciones MAP para  $\theta$  se deben encontrar los valores de los parámetros para los cuales se cumpla que

$$\frac{\partial \log \mathbf{G}(\theta)}{\partial \theta} = \frac{\partial \log \mathbf{L}(\theta; \mathbf{y}, \mathbf{z})}{\partial \theta} + \frac{\partial \log p(\theta)}{\partial \theta} = 0\tag{4.5}$$

Aquí

$$\begin{aligned}
\frac{\partial \log \mathbf{L}(\boldsymbol{\theta}; \mathbf{y}, \mathbf{z})}{\partial \boldsymbol{\theta}} &= \sum_{i=1}^N \frac{\partial}{\partial \boldsymbol{\theta}} (\log P(\mathbf{Y}_i = \mathbf{y}_i | Z_i = z_i)) \\
&= \sum_{i=1}^N \frac{1}{P(\mathbf{Y}_i = \mathbf{y}_i | Z_i = z_i)} \frac{\partial}{\partial \boldsymbol{\theta}} P(\mathbf{Y}_i = \mathbf{y}_i | Z_i = z_i) \\
&= \sum_{i=1}^N \frac{1}{P(\mathbf{Y}_i = \mathbf{y}_i | Z_i = z_i)} \frac{\partial}{\partial \boldsymbol{\theta}} \left( \sum_{c=1}^C p(\boldsymbol{\alpha}_c | Z_i = z_i) P(\mathbf{Y}_i = \mathbf{y}_i | \boldsymbol{\alpha}_c, Z_i = z_i) \right) \\
&= \sum_{i=1}^N \frac{1}{P(\mathbf{Y}_i = \mathbf{y}_i | Z_i = z_i)} \sum_{c=1}^C \frac{\partial}{\partial \boldsymbol{\theta}} p(\boldsymbol{\alpha}_c | Z_i = z_i) P(\mathbf{Y}_i = \mathbf{y}_i | \boldsymbol{\alpha}_c, Z_i = z_i)
\end{aligned}$$

y la probabilidad a posteriori de pertenencia a las clases que la denotamos por  $P_{ic}$ , satisface

$$P_{ic} = p(\boldsymbol{\alpha}_c | Z_i = z_i, \mathbf{Y}_i = \mathbf{y}_i) = \frac{p(\boldsymbol{\alpha}_c | Z_i = z_i) P(\mathbf{Y}_i = \mathbf{y}_i | \boldsymbol{\alpha}_c, Z_i = z_i)}{P(\mathbf{Y}_i = \mathbf{y}_i | Z_i = z_i)} \quad (4.6)$$

El algoritmo EM es un método general para obtener estimaciones MV con datos faltantes (Dempster et al.; 1977; McLachlan y Krishnan; 1997). Dado que la membresía de la clase es un dato no observable, el modelo (3.9), que es una analogía al de extensión de regresión del análisis de clases latentes con predictores (LCRP), se convierte en un problema típico de datos incompletos. Por tanto, el algoritmo de Esperanza-Maximización o Esperanza-Minimización (EM) es un enfoque iterativo ideal para calcular las estimaciones de MAP o MV. El algoritmo EM maximiza la probabilidad al iterar entre la imputación de datos faltantes de un modelo parametrizado en las estimaciones más recientes y la maximización de la probabilidad de datos completos (conjunta con respecto a los datos observables y faltantes). Formalmente, la imputación se realiza a través de un paso E (esperanza) que calcula la probabilidad esperada de los datos completos dados los datos observados y un paso M (maximización) que maximiza la probabilidad calculada a partir del paso E.

Para la estimación EM se introduce una variable aleatoria dicotómica  $S_{ic}$  que nos indique si el individuo  $i$  pertenece (1) o no (0) a la clase latente  $c$  y sea  $\boldsymbol{\theta}$  el vector de parámetros descrito en (4.2). Aquí los  $S_{ic}$ , cuyos valores observables los denotaremos por  $s_{ic}$  se consideran como “datos faltantes”. Por tanto, si  $S_{ic}$  fuera directamente observable, la distribución a posteriori completa del modelo RDINA con covariable sería proporcional a

$$\begin{aligned}
\mathbf{G}_c(\boldsymbol{\theta}) &= \mathbf{L}_c(\boldsymbol{\theta}; \mathbf{y}, \mathbf{z}, \mathbf{s}) p(\boldsymbol{\theta}) \\
&= \prod_{i=1}^N \prod_{c=1}^C \left( p(\boldsymbol{\alpha}_c | Z_i = z_i) \prod_{j=1}^J P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_c, Z_i = z_i) \right)^{S_{ic}} p(\boldsymbol{\theta}), \quad (4.7)
\end{aligned}$$

y su logaritmo quedaría representada por :

$$\begin{aligned}
\log(\mathbf{G}_c(\boldsymbol{\theta})) &= \sum_{i=1}^N \sum_{c=1}^C \left( S_{ic} \log p(\boldsymbol{\alpha}_c | Z_i = z_i) + \sum_{j=1}^J S_{ic} \log P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_c, Z_i = z_i) \right) \\
&+ \log p(\boldsymbol{\theta}) \\
&= \sum_{i=1}^N \sum_{c=1}^C S_{ic} \log p(\boldsymbol{\alpha}_c | Z_i = z_i) \\
&+ \sum_{i=1}^N \sum_{c=1}^C \sum_{j=1}^J S_{ic} \log P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_c, Z_i = z_i) + \log p(\boldsymbol{\theta}), \quad (4.8)
\end{aligned}$$

Introduciéndose la función de enlace *logit* dentro del modelo, la función de log verosimilitud completa queda más explícitamente caracterizada salvo una constante por:

$$\begin{aligned} \log(\mathbf{G}_c(\boldsymbol{\theta})) &= \sum_{i=1}^N \sum_{c=1}^C \left( s_{ic} \sum_{k=1}^K \log \left( \frac{e^{b_k + h_k z_i}}{1 + e^{b_k + h_k z_i}} \right) \right) + \\ &\sum_{i=1}^N \sum_{c=1}^C \sum_{j=1}^J \left( s_{ic} \log \left( \left( \frac{e^{f_j + d_j \eta_{cj} + \ell_j z_i}}{1 + e^{f_j + d_j \eta_{cj} + \ell_j z_i}} \right)^{y_{ij}} \left( \frac{1}{1 + e^{f_j + d_j \eta_{cj} + \ell_j z_i}} \right)^{1 - y_{ij}} \right) \right) + \log p(\boldsymbol{\theta}). \end{aligned} \quad (4.9)$$

En algunas ocasiones estas ecuaciones pueden resolverse analíticamente o numéricamente, pero cuando la solución no es posible hay dos enfoques principales para obtenerla: el método Newton y el algoritmo EM, el cual aplicaremos a continuación.

### Paso Esperanza

La idea básica del algoritmo EM es considerar las observaciones compuestas por datos observables  $\mathbf{y}$ , así como también por datos faltantes  $\mathbf{s}$ . En algunos casos  $\mathbf{s}$  podría tratarse de datos que realmente faltan, pero en otros son solo datos adicionales (aumentados) que desearíamos tener. El primer paso es determinar la posteriori conjunta de los datos completos ( $\mathbf{y}, \mathbf{s}$ ) que es proporcional a  $\mathbf{G}_c(\boldsymbol{\theta})$  en (4.7) si se conociera  $\mathbf{s}$ , dado que no se tiene, en realidad no puede calcularse  $\mathbf{G}_c(\boldsymbol{\theta})$ . Se necesita siempre maximizar una función que dependa solo de  $\boldsymbol{\theta}$  y de los datos observados  $\mathbf{y}$ . Por lo tanto, el paso **E** del algoritmo EM consiste en calcular la esperanza condicional de  $\log(\mathbf{G}_c(\boldsymbol{\theta}))$  dado  $\mathbf{y}$  suponiendo que el verdadero valor del parámetro es conocido en la iteración  $p$  del algoritmo. Definamos entonces

$$Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(p)}) = E_{\boldsymbol{\theta}^{(p)}} [\log(\mathbf{G}_c(\boldsymbol{\theta})) | \mathbf{Y} = \mathbf{y}, \mathbf{Z} = \mathbf{z}], \quad (4.10)$$

donde el operador de la esperanza **E** tiene el subíndice  $\boldsymbol{\theta}^{(p)}$  para señalar explícitamente que esta esperanza es evaluada en  $\boldsymbol{\theta}^{(p)}$  como  $\boldsymbol{\theta}$ . A partir de aquí, se deduce que en la  $(p+1)$ -ésima iteración, el paso **E** requiere el cálculo de  $Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(p)})$ , donde  $\boldsymbol{\theta}^{(p)}$  es el valor de  $\boldsymbol{\theta}$  después de la  $p$ -ésima iteración EM. Como la log-posteriori del conjunto de datos completos,  $\log(\mathbf{G}_c(\boldsymbol{\theta}))$  es lineal en la variable no observable  $S_{ic}$ , el paso E en la  $(p+1)$ -ésima iteración, simplemente requiere del cálculo de la actual esperanza condicional de  $S_{ic}$  dados los datos observados  $\mathbf{y}$  y  $\mathbf{z}$ , el cual se determina usando las ecuaciones (4.10) y (4.8) de la siguiente manera:

$$\begin{aligned} Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(p)}) &= E_{\boldsymbol{\theta}^{(p)}} \left( \sum_{i=1}^N \sum_{c=1}^C S_{ic} \log p(\boldsymbol{\alpha}_c | Z_i = z_i) | \mathbf{Y} = \mathbf{y}, \mathbf{Z} = \mathbf{z} \right) \\ &+ E_{\boldsymbol{\theta}^{(p)}} \left( \sum_{i=1}^N \sum_{c=1}^C \sum_{j=1}^J S_{ic} \log P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_c, Z_i = z_i) | \mathbf{Y} = \mathbf{y}, \mathbf{Z} = \mathbf{z} \right) + \log p(\boldsymbol{\theta}) \\ &= \sum_{i=1}^N \sum_{c=1}^C E_{\boldsymbol{\theta}^{(p)}}(S_{ic} | Z_i = z_i, \mathbf{Y}_i = \mathbf{y}_i) \log p(\boldsymbol{\alpha}_c | Z_i = z_i) \\ &+ \sum_{i=1}^N \sum_{c=1}^C \sum_{j=1}^J E_{\boldsymbol{\theta}^{(p)}}(S_{ic} | Z_i = z_i, \mathbf{Y}_i = \mathbf{y}_i) \log P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_c, Z_i = z_i) + \log p(\boldsymbol{\theta}) \\ &= \sum_{i=1}^N \sum_{c=1}^C P_{ic} \log p(\boldsymbol{\alpha}_c | Z_i = z) \\ &+ \sum_{i=1}^N \sum_{c=1}^C \sum_{j=1}^J P_{ic} \log P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_c, Z_i = z) + \log p(\boldsymbol{\theta}) \end{aligned} \quad (4.11)$$

donde  $P_{ic} = E_{\theta^{(p)}}(S_{ic}|Z_i = z_i, \mathbf{Y}_i = \mathbf{y}_i)$  se calcula por (4.6) en base a las estimaciones.

### Paso Maximización

La  $(p + 1)$ -ésima iteración del paso M requiere de la maximización global de  $Q(\theta|\theta^{(p)})$  con respecto a  $\theta$  para obtener la actualización del estimado  $\theta^{(p+1)}$ .

Debido a que no existe, en este problema, una solución completa para el proceso de maximización anterior, utilizaremos el método de Newton-Raphson de una iteración (Lange; 1995) para aproximar los valores máximos en el paso M.

Huang y Bandeen-Roche (2004) demostraron que esta aproximación de un solo paso tiene una tasa de convergencia que es casi idéntica a la tasa del algoritmo EM y, por lo tanto, ahorra tiempo al realizar repetidamente el método de Newton. Los pasos E y M se alternan repetidamente hasta que la diferencia

$$\mathbf{G}(\theta^{(p+1)}) - \mathbf{G}(\theta^{(p)})$$

sea una cantidad arbitrariamente pequeña. Dempster et al. (1977) demostró que la función verosimilitud de datos incompletos  $\mathbf{G}(\theta)$  no disminuye después de una iteración EM; esto es:

$$\mathbf{G}(\theta^{(p+1)}) \geq \mathbf{G}(\theta^{(p)}) \quad (4.12)$$

para  $p = 0, 1, 2, \dots$ . Por lo tanto, la convergencia debe obtenerse con una secuencia de valores de la verosimilitud que están superiormente acotados.

De esta manera, el paso de Maximización (M) implica encontrar un nuevo  $\theta^{(p+1)}$  que maximice  $Q(\theta|\theta^{(p)})$ . A veces, se dispone de soluciones de forma cerrada en el paso M. En otros casos, se pueden usar métodos iterativos estándar para mejorar la log-posterior de datos completos. El software Latent GOLD utiliza ajuste proporcional iterativo (IPF) y Newton unidimensional en el paso M (Vermunt; 1997).

Además del algoritmo EM, el software también utiliza el método de Newton-Raphson (NR), que permite la actualización de los parámetros de la siguiente manera:

$$\theta^{(p+1)} = \theta^{(p)} - \varepsilon \mathbf{H}^{-1} \boldsymbol{\nu} \quad (4.13)$$

donde el vector de gradiente  $\boldsymbol{\nu}$  contiene las derivadas de primer orden del log posterior a todos los parámetros evaluados en  $\theta^{(p)}$ ,  $\mathbf{H}$  es la matriz Hessian que contiene las derivadas de segundo orden a todos los parámetros y  $\varepsilon$  es un escalar que indica el tamaño del paso. Cada elemento  $\nu_c$  de  $\boldsymbol{\nu}$  es igual a

$$\nu_c = \sum_{i=1}^N \frac{\partial \log P(\mathbf{Y}_i = \mathbf{y}_i | Z_i = z)}{\partial \theta_c} + \frac{\partial \log p(\theta)}{\partial \theta_c} \quad (4.14)$$

y el elemento  $H_{cc'}$  de  $\mathbf{H}$  es igual

$$H_{cc'} = \sum_{i=1}^N \frac{\partial^2 \log P(\mathbf{Y}_i = \mathbf{y}_i | Z_i = z)}{\partial \theta_c \partial \theta_{c'}} + \frac{\partial^2 \log p(\theta)}{\partial \theta_c \partial \theta_{c'}} \quad (4.15)$$

El software Latent Gold calcula estas derivadas analíticamente. El tamaño del paso  $\varepsilon$  ( $0 < \varepsilon \leq 1$ ) es necesario para evitar que ocurran disminuciones de la log a posteriori, es decir, cuando debido a una actualización NR estándar,  $\mathbf{H}^{-1} \boldsymbol{\nu}$  produce una disminución de la log verosimilitud, el tamaño del paso se reduce hasta que esto ya no ocurra.



La matriz  $-\mathbf{H}^{-1}$  evaluada en el  $\hat{\boldsymbol{\theta}}$  final nos da la estimación estándar para la matriz de varianza covarianza asintótica de los parámetros del modelo:  $\hat{\Sigma}_{standard}(\boldsymbol{\theta}) = -\hat{\mathbf{H}}^{-1}$ .

## Distribuciones priori

Con el fin de evitar soluciones límite, es decir, el problema de la inexistencia de estimaciones MV, el software Latent Gold ha implementado algunas ideas de estadística Bayesiana. Los problemas de límites se refieren a cuando las probabilidades multinomiales pueden convertirse en cero, para lo cual el software usa una priori de Dirichlet. Esta también es llamada una priori conjugada, ya que tienen la misma forma que la densidad de probabilidad multinomial. La implicación del uso de una priori hace que el método de estimación ya no sea MV, sino de la moda a posteriori (MAP).

La influencia de la priori en las estimaciones finales de los parámetros depende de los valores elegidos para los hiperparámetros  $\alpha$ 's de la distribución Dirichlet, así como el tamaño de la muestra. La configuración predeterminada es una simétrica con  $\alpha = 1$ . Esto demuestra que, con tamaños de muestra moderados, la influencia de las priors en la estimación de los parámetros es insignificante. El ajuste  $\alpha = 0$  produce estimaciones de MV.

## Convergencia

El programa cuyo algoritmo se ha implementado en Latent Gold comienza con EM hasta alcanzar el número máximo de iteraciones EM o el criterio de convergencia EM (Tolerancia). Luego, el programa cambia a iteraciones NR que se detienen cuando se alcanza el número máximo de iteraciones NR o el criterio de convergencia general (Tolerancia).

El criterio de convergencia que se utiliza es

$$\sum_{r=1}^{npar} \left| \frac{\theta_r^{(p)} - \theta_r^{(p-1)}}{\theta_r^{(p-1)}} \right| \quad (4.16)$$

donde  $npar$  es el número de parámetros estimados. La ecuación (4.16) es la suma de los cambios relativos absolutos en los parámetros. Cabe señalar que a veces es más eficiente usar solo el algoritmo EM, que se logra estableciendo como límites de iteración Newton-Raphson = 0. Este caso podría aplicarse en modelos con muchos parámetros.

## Valores iniciales

El software Latent Gold usa la opción técnica *Seed* igual a 0 para generar valores iniciales aleatorios que difieren cada vez que se estima un modelo. Dado que el algoritmo EM es estable, el uso de valores iniciales generalmente es lo suficiente bueno como para obtener una solución convergente.

La mejor manera de evitar terminar con una solución local es el uso de múltiples conjuntos de valores iniciales, pues diferentes conjuntos de valores iniciales pueden generar soluciones con diferentes valores log posterioris. En Latent Gold, el uso de tales conjuntos múltiples de valores iniciales aleatorios está automatizado. Puede especificarse cuántos conjuntos de valores iniciales debe usar el programa en la opción *Random sets*.

Otro parámetro relevante es *Iterations* que especifica el número de iteraciones que se realizarán por conjunto de inicio. Dentro de cada uno de los conjuntos aleatorios, Latent GOLD realiza el número especificado de iteraciones EM. Posteriormente, con el mejor 10 por ciento (redondeado hacia arriba) en términos de la log posteriori, el programa realiza 2 veces más *Iterations* mediante iteraciones EM. Finalmente, continúa con la mejor solución hasta

la convergencia. Este procedimiento aumenta considerablemente la probabilidad de encontrar la solución global de MAP o MV, especialmente si ambos parámetros se configuran lo suficientemente grandes, pero en general no garantiza que se encontrará en una sola ejecución.

Con la opción *Tolerance* se puede especificar el criterio de convergencia EM que se utilizará dentro del procedimiento de valores iniciales aleatorios. Por lo tanto, las iteraciones de valores iniciales se detienen si se alcanza esta tolerancia o el número máximo de iteraciones.

#### 4.3.3. Estimación de la varianza

Debido a que el algoritmo EM es un método para la estimación por máxima verosimilitud o MAP en problemas de datos incompletos, la matriz de información de Fisher observada basada en la función de verosimilitud  $\mathbf{G}(\boldsymbol{\theta})$  de datos incompletos puede usarse para estimar errores estándar de las estimaciones de los parámetros que condicionan el número de clases. Sin embargo, la evaluación de las derivadas de segundo orden de la función de *log* verosimilitud de datos incompletos puede ser difícil de obtener. Louis (1982) demostró que la matriz de información observada de datos incompletos puede calcularse en términos de las derivadas parciales de primer y segundo orden de la función de *log* verosimilitud de datos completos que fue introducido en el marco EM. Por lo tanto, se implementará el enfoque de Louis para calcular la matriz de varianza-covarianza de las estimaciones de los parámetros.

Para tal implementación denotamos por  $D_{\boldsymbol{\theta}}^1$  y  $D_{\boldsymbol{\theta}}^2$  al gradiente y al operador Hessiano con respecto a  $\boldsymbol{\theta}$ , y definamos  $\mathbf{I}(\boldsymbol{\theta}; \mathbf{y}, \mathbf{z}) = -D_{\boldsymbol{\theta}}^2 \log \mathbf{G}(\boldsymbol{\theta})$  la matriz de información de Fisher observada de los datos incompletos de la función de verosimilitud con respecto a los elementos de  $\boldsymbol{\theta}$ , y la de datos completos  $\mathbf{I}_c(\boldsymbol{\theta}; \mathbf{y}, \mathbf{z}, \mathbf{s}) = -D_{\boldsymbol{\theta}}^2 \log \mathbf{G}_c(\boldsymbol{\theta})$ .

Louis (1982) demostró que

$$\mathbf{I}(\hat{\boldsymbol{\theta}}; \mathbf{y}, \mathbf{z}) = \mathcal{I}_c(\hat{\boldsymbol{\theta}}; \mathbf{y}, \mathbf{z}) - \text{Var}(\mathbf{s}_c(\hat{\boldsymbol{\theta}}; \mathbf{y}, \mathbf{z}, \mathbf{s}) | \mathbf{y}) ,$$

donde

$$\begin{aligned} \mathcal{I}_c(\hat{\boldsymbol{\theta}}; \mathbf{y}, \mathbf{z}) &= E(\mathbf{I}_c(\boldsymbol{\theta}; \mathbf{y}, \mathbf{z}, \mathbf{s}) | \mathbf{y})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} , \\ \mathbf{s}_c(\hat{\boldsymbol{\theta}}; \mathbf{y}, \mathbf{z}, \mathbf{s}) &= D_{\boldsymbol{\theta}}^1 \log \mathbf{G}_c(\boldsymbol{\theta})|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}} \end{aligned}$$

y  $\hat{\boldsymbol{\theta}}$  es el MAP de  $\boldsymbol{\theta}$

Por lo tanto, la matriz de información observada se puede calcular en términos de las derivadas parciales de primer y segundo orden de la función de verosimilitud de datos completos introducida dentro del marco del algoritmo EM. Louis (1982) muestra además que  $\mathcal{I}_c(\hat{\boldsymbol{\theta}}; \mathbf{y}, \mathbf{z})$  tiene la misma fórmula que la segunda derivada de  $Q(\boldsymbol{\theta} | \hat{\boldsymbol{\theta}})$ .

Considerando la teoría asintótica de los estimadores de máxima verosimilitud, el inverso de la matriz de información observada evaluado en las estimaciones de los parámetros,  $(-D_{\boldsymbol{\theta}}^2 \log \mathbf{G}(\hat{\boldsymbol{\theta}}))^{-1}$ , vendría a ser el estimador para la matriz de varianza-covarianza del estimador de la maximización a posteriori  $\hat{\boldsymbol{\theta}}$ , con lo cual los errores estándar de los parámetros de este estimador se pueden hallar como las raíces cuadradas de las entradas en la diagonal de esta matriz.

#### 4.4. Selección del modelo

Entre las medidas más usadas para encontrar un modelo que tenga el mejor balance entre un conjunto de datos de tamaño  $n$ , el ajuste del modelo y el número de parámetros a ser estimados, se encuentran el Bayesian Information Criterion (BIC) y el Akaike Information Criterion (AIC).

$$AIC = D(\hat{\boldsymbol{\theta}}) + 2\Upsilon , \quad (4.17)$$

$$BIC = D(\hat{\boldsymbol{\theta}}) + \Upsilon \log(n) , \quad (4.18)$$

donde  $D(\hat{\boldsymbol{\theta}}) = -2\log[\mathbf{G}(\hat{\boldsymbol{\theta}})]$ , siendo  $\mathbf{G}(\hat{\boldsymbol{\theta}})$  la función de verosimilitud del modelo y  $\hat{\boldsymbol{\theta}}$  es el estimador de la moda a posteriori. Asimismo,  $\Upsilon$  es el número de parámetros del modelo.

Ambos indicadores son medidas de calidad del modelo estadístico, que buscan un equilibrio entre la bondad de ajuste del modelo y su complejidad, penalizando la cantidad de parámetros calculados y el volumen de datos. Si bien no se tratan de pruebas estadísticas, estos indicadores nos señalan cuál es el mejor de todos los modelos analizados, el cual será aquel con el menor valor de AIC y BIC, es decir, una disminución en el valor de  $D(\hat{\boldsymbol{\theta}})$  producirá un mejor ajuste en el modelo. Adicionalmente, para el caso del BIC, se tiene una penalización por el incremento de parámetros que incluso es mayor, dado que se requiere controlar el sobreajuste.



## Capítulo 5

# Estudio de Simulación

En el presente capítulo se desarrollan estudios de simulación de la extensión del modelo RDINA con una covariable afectando habilidades e ítems, con la finalidad de estudiar la precisión en la recuperación del vector de parámetros de los ítems ( $\mathbf{f}$ ,  $\mathbf{d}$ ,  $\ell$ ) y los parámetros de los atributos ( $\mathbf{b}$ ,  $\mathbf{h}$ ).

Para la generación de datos se usan como valores iniciales los parámetros obtenidos por [Park y Lee \(2014\)](#) en su ejercicio de aplicación de la extensión del modelo DINA con covariable. Este está basado en el análisis de datos para alumnos del cuarto grado en el Estudio de Tendencias Internacionales de Matemáticas y Ciencias (TIMMS, 2007).

Basándonos en ese contexto se usa lo siguiente:

- Número de habilidades: 7.
- Número de ítems: Se usará un solo tamaño correspondiente a  $J = 25$  ítems.
- Número de réplicas: [Harwell et al. \(1996\)](#) recomendaron realizar como mínimo 25 réplicas para obtener resultados confiables en la TRI. Por otro lado [Ma et al. \(2016\)](#), realizan 100 réplicas en su estudio de simulación por cada condición con la finalidad de estimar los parámetros usando diferentes modelos de diagnóstico cognitivo. En este estudio se usan también 100 réplicas.
- Matriz  $Q$ : matriz de orden  $25 \times 7$  basado en [Lee et al. \(2011\)](#) para el TIMMS del cuarto grado en el área de Matemáticas.

En primer lugar, se simulan matrices de respuestas dicotómicas a fin de evaluar la precisión en la recuperación de los parámetros bajo los escenarios siguientes, usando el software libre R:

- a) El modelo RDINA con covariable afectando a los ítems y atributos
- b) El modelo RDINA con covariable afectando solo a los ítems
- c) El modelo RDINA con covariable afectando solo a los atributos
- d) El modelo RDINA sin covariable.

Luego, para la estimación se empleará el software Latent GOLD 5.1 que utilizan los algoritmos de Esperanza-Maximización (EM) y de Newton-Raphson descritos en la sección anterior en la obtención de estimaciones mediante el método de la moda a posteriori (MAP). El uso de la estimación MAP evita los problemas de límites comúnmente asociados con los modelos de clase latentes (DeCarlo; 2011).

## 5.1. Algoritmo para simular datos

Para poder simular la matriz de respuestas dicotómicas que se usa en el modelo RDINA con covariable se sigue el siguiente procedimiento:

- Definir la cantidad de ítems  $J$ , el número de habilidades  $K$  y la matriz  $Q$ .
- Fijar los parámetros iniciales de los ítems ( $\mathbf{f}, \mathbf{d}, \ell$ ) y de las habilidades ( $\mathbf{b}, \mathbf{h}$ ) en los valores reales obtenidos por Yoon Soo Park (2014) en su aplicación sobre el TIMMS para cada uno de los 3 escenarios.
- Generar una sola vez una covariable continua  $\mathbf{Z}$  simulando valores de una distribución normal  $N(17.18, 25)$  para los  $N$  examinados.
- Calcular, para cada uno de los  $N$  examinados, las probabilidades de que posean la habilidad  $k$  ( $p(\alpha_{ck}|Z_i = z)$ ) dada su pertenencia a la clase  $c$ , usando los parámetros de las habilidades ( $\mathbf{b}, \mathbf{h}$ ), el valor de la covariable inicial  $\mathbf{Z}$  y la ecuación (3.11).
- Calcular para cada uno de los  $N$  examinados la variable latente binaria  $\eta_{cj}$  simulando ensayos de Bernoulli a las probabilidades de habilidades calculado en el punto anterior y considerando la matriz  $Q$ .
- Calcular las probabilidades de respuesta  $P(Y_{ij} = 1|\alpha_c, Z_i = z)$  usando (3.9).
- Simular para cada uno de los  $N$  examinados las respuestas a cada uno de los  $J$  ítems mediante ensayos de Bernoulli aplicados a las probabilidades calculadas en el punto anterior.

Los códigos para simular la matriz de respuestas dicotómicas para el modelo RDINA con covariable, así como los códigos para la generación de los valores iniciales se muestran en el Apéndice B.

## 5.2. Criterios para evaluar la simulación

Con la finalidad de estudiar la precisión en la estimación de los parámetros obtenidos, se calcula el error cuadrático medio (MSE) y el sesgo. El MSE se define como la sumatoria de las diferencias al cuadrado entre los parámetros estimados y los verdaderos divididos entre la cantidad de réplicas.

$$MSE = \sum_{r=1}^R \frac{(\hat{\theta}_r - \hat{\theta})^2}{R} \quad (5.1)$$

donde  $\theta$  representa a los parámetros de los ítems y de los atributos y  $\hat{\theta}$  es el estimado del vector de parámetros anterior, siendo  $R$  el número de réplicas a usarse en la simulación. Este indicador se usa para evaluar la varianza total del error de estimación entre los valores observados y los reales. De esta manera, resultados con valor cero nos da un indicativo de que no hay discrepancia entre los valores verdaderos y el modelo.

El sesgo se define como la sumatoria de las diferencias entre el valor estimado y el verdadero dividido entre la cantidad de réplicas.

$$Sesgo = \sum_{r=1}^R \frac{\hat{\theta}_r - \hat{\theta}}{R} \quad (5.2)$$

Y el sesgo porcentual se define como el valor absoluto que resulta de dividir el sesgo entre el valor verdadero.

$$\%Sesgo = \left| \frac{Sesgo}{valor\ real} \right| * 100 \quad (5.3)$$

### 5.3. Simulaciones

A continuación, se muestran los resultados de las simulaciones para la recuperación de cada uno de los parámetros en el caso que  $N = 500$ ,  $K = 7$  y  $J = 25$ , y en la parte inferior de cada tabla se ha calculado  $\mathbf{b_K} = \sum_{k=1}^K b_k/K$ ,  $\mathbf{h_K} = \sum_{k=1}^K h_k/K$ ,  $\mathbf{f_J} = \sum_{j=1}^J f_j/J$ ,  $\mathbf{d_J} = \sum_{j=1}^J d_j/J$ ,  $\mathbf{\ell_J} = \sum_{j=1}^J \ell_j/J$ , que representa el promedio de los parámetros contenidos en los vectores  $\mathbf{b}$ ,  $\mathbf{h}$ ,  $\mathbf{f}$ ,  $\mathbf{d}$  y  $\mathbf{\ell}$  respectivamente y nos sirve, además, para resumir los resultados generales de dicho vector.

#### 5.3.1. Modelo RDINA con covariable afectando ítems y atributos

La tabla 5.1 presenta los resultados de la precisión en la recuperación de parámetros del vector  $\mathbf{b}$  y muestra que el sesgo porcentual de los parámetros fue igual o inferior al 4.8 %, excepto en el parámetro  $b_6$  que fue de 16.6 %. Además, el sesgo porcentual de  $\mathbf{b_K}$  que representa al promedio de las estimaciones de los parámetros fue de 2.57 %.

**Tabla 5.1:** Estimación del vector de parámetros  $\mathbf{b}$  del modelo RDINA con covariable afectando ítems y atributos en la simulación

Parámetro	Valor real	MSE	Sesgo	% Sesgo
$b_1$	-5.65	1.077	0.088	1.6 %
$b_2$	-5.52	5.221	-0.122	2.2 %
$b_3$	-3.06	4.652	0.131	4.3 %
$b_4$	-0.55	5.91	0.022	3.9 %
$b_5$	-3.04	3.051	0.146	4.8 %
$b_6$	-2.77	12.317	0.46	16.6 %
$b_7$	-3.58	3.728	-0.103	2.9 %
$\mathbf{b_K}$	<b>-3.4529</b>	<b>0.5877</b>	<b>0.0888</b>	<b>2.57 %</b>

La tabla 5.2 presenta la recuperación de parámetros del vector  $\mathbf{h}$  y muestra que el sesgo porcentual de los parámetros fue igual o inferior al 3.3 %. Además, el sesgo porcentual de  $\mathbf{h_K}$  que representa al promedio de las estimaciones de los parámetros fue de 1.23 %.

**Tabla 5.2:** Estimación del vector de parámetros  $\mathbf{h}$  del modelo RDINA con covariable afectando ítems y atributos en la simulación

Parámetro	Valor real	MSE	Sesgo	% Sesgo
$h_1$	0.34	0.003	-0.005	1.4%
$h_2$	0.33	0.015	0	0.1%
$h_3$	0.24	0.015	-0.006	2.5%
$h_4$	0.08	0.032	-0.001	1.2%
$h_5$	0.21	0.01	-0.008	4%
$h_6$	0.18	0.045	-0.006	3.3%
$h_7$	0.23	0.011	0.006	2.6%
$\mathbf{h_K}$	<b>0.23</b>	<b>0.0025</b>	<b>-0.0028</b>	<b>1.23%</b>

La tabla 5.3 muestra la recuperación de parámetros del vector  $\mathbf{f}$  y se aprecia que el sesgo porcentual de los parámetros fue igual o inferior al 6.3 %. Además, el sesgo porcentual de  $\mathbf{f_J}$  que representa al promedio de las estimaciones de los parámetros fue de 2.33 %.

**Tabla 5.3:** Estimación del vector de parámetros  $\mathbf{f}$  del modelo RDINA con covariable afectando ítems y atributos en la simulación

Parámetro	Valor real	MSE	Sesgo	% Sesgo
$f_1$	-1.58	0.221	-0.008	0.5%
$f_2$	-5.82	1.086	-0.21	3.6%
$f_3$	-3.32	0.4	-0.108	3.2%
$f_4$	-1.77	0.169	-0.047	2.7%
$f_5$	-4.36	0.542	-0.146	3.3%
$f_6$	-2.97	1.101	-0.12	4.1%
$f_7$	-3.19	0.392	-0.141	4.4%
$f_8$	-3.17	0.246	-0.139	4.4%
$f_9$	-3.58	1.567	-0.208	5.8%
$f_{10}$	-2.41	0.501	-0.059	2.4%
$f_{11}$	-1.63	0.116	-0.084	5.2%
$f_{12}$	-1.01	0.178	0.035	3.4%
$f_{13}$	-2.36	0.221	-0.064	2.7%
$f_{14}$	-3.06	0.262	-0.073	2.4%
$f_{15}$	-2.36	0.25	-0.062	2.6%
$f_{16}$	-5.89	0.624	0.023	0.4%
$f_{17}$	-4.66	0.37	0.011	0.2%
$f_{18}$	-3.45	0.325	-0.004	0.1%
$f_{19}$	-5.11	0.438	0.057	1.1%
$f_{20}$	-1.94	0.216	0.038	1.9%
$f_{21}$	-5.02	0.419	-0.157	3.1%
$f_{22}$	-1.74	0.16	-0.11	6.3%
$f_{23}$	-4.58	0.383	-0.039	0.9%
$f_{24}$	-3.36	1.105	-0.207	6.2%
$f_{25}$	-3.27	0.285	-0.081	2.5%
$\mathbf{f_J}$	<b>-3.2644</b>	<b>0.0324</b>	<b>-0.0761</b>	<b>2.33%</b>

La tabla 5.4 muestra los resultados en la recuperación de parámetros del vector  $\mathbf{d}$  y se aprecia que el sesgo porcentual de los parámetros fue igual o inferior al 5.5 %, excepto los parámetros  $d_{10}$ ,  $d_7$ ,  $d_{25}$  y  $d_{20}$  que fueron de 13.3 %, 13.2 %, 9 % y 6.9 %, respectivamente. Además, el sesgo porcentual de  $\mathbf{d_J}$  que representa al promedio de las estimaciones de los parámetros fue de 3.86 %.

La tabla 5.5 presenta la recuperación de parámetros del vector  $\ell$  y muestra que el sesgo porcentual de los parámetros fue igual o inferior al 6.3 %, excepto los parámetros  $\ell_{22}$  y  $\ell_{11}$  que fueron de 8.7 % y 7.4 %, respectivamente. Además, el sesgo porcentual de  $\ell_J$  que representa al promedio de las estimaciones de los parámetros fue de 2.48 %.

**Tabla 5.4:** Estimación del vector de parámetros  $\mathbf{d}$  del modelo RDINA con covariable afectando ítems y atributos en la simulación

Parámetro	Valor real	MSE	Sesgo	% Sesgo
$d_1$	1.51	0.18	-0.009	0.6 %
$d_2$	2.26	0.96	0.096	4.2 %
$d_3$	2.06	0.252	0.075	3.7 %
$d_4$	1.44	0.441	0.058	4 %
$d_5$	2.63	0.23	0.095	3.6 %
$d_6$	2.87	4.816	-0.072	2.5 %
$d_7$	2.54	3.048	0.337	13.2 %
$d_8$	1.73	0.212	0.04	2.3 %
$d_9$	2.66	1.805	0.1	3.8 %
$d_{10}$	1.56	1.093	0.208	13.3 %
$d_{11}$	0.45	0.179	0.01	2.3 %
$d_{12}$	1.04	0.209	0.051	4.9 %
$d_{13}$	2.05	0.575	0.044	2.1 %
$d_{14}$	1.55	0.173	0.024	1.6 %
$d_{15}$	1.56	0.132	0.041	2.6 %
$d_{16}$	2.53	0.17	-0.008	0.3 %
$d_{17}$	1.76	0.127	0.002	0.1 %
$d_{18}$	1.86	0.21	0.055	3 %
$d_{19}$	2.19	0.213	0.113	5.2 %
$d_{20}$	1.83	0.254	0.126	6.9 %
$d_{21}$	1.08	0.206	0.037	3.4 %
$d_{22}$	0.77	0.367	-0.021	2.7 %
$d_{23}$	2.31	0.239	0.057	2.5 %
$d_{24}$	1.88	0.891	0.103	5.5 %
$d_{25}$	2.81	1.059	0.252	9 %
$\mathbf{d_J}$	<b>1.8772</b>	<b>0.0498</b>	<b>0.0725</b>	<b>3.86 %</b>

**Tabla 5.5:** Estimación del vector de parámetros  $\ell$  del modelo RDINA con covariable afectando ítems y atributos en la simulación

Parámetro	Valor real	MSE	Sesgo	% Sesgo
$\ell_1$	0.15	0.001	0	0.3 %
$\ell_2$	0.21	0.003	0.01	4.9 %
$\ell_3$	0.13	0.001	0.004	3.3 %
$\ell_4$	0.13	0.001	0.005	3.5 %
$\ell_5$	0.23	0.002	0.008	3.3 %
$\ell_6$	0.31	0.004	0.013	4.2 %
$\ell_7$	0.17	0.002	0.01	5.8 %
$\ell_8$	0.15	0.001	0.006	4 %
$\ell_9$	0.21	0.004	0.013	6.3 %
$\ell_{10}$	0.12	0.002	0.004	2.9 %
$\ell_{11}$	0.05	0.001	0.004	7.4 %
$\ell_{12}$	0.12	0.001	-0.002	1.4 %
$\ell_{13}$	0.18	0.001	0.005	2.8 %
$\ell_{14}$	0.13	0.001	0.006	4.3 %
$\ell_{15}$	0.17	0.001	0.004	2.1 %
$\ell_{16}$	0.28	0.002	-0.002	0.6 %
$\ell_{17}$	0.15	0.001	0	0 %
$\ell_{18}$	0.19	0.001	0	0.1 %
$\ell_{19}$	0.23	0.002	-0.005	2 %
$\ell_{20}$	0.11	0.001	-0.003	2.9 %
$\ell_{21}$	0.16	0.001	0.006	3.6 %
$\ell_{22}$	0.08	0.001	0.007	8.7 %
$\ell_{23}$	0.24	0.002	0.002	0.9 %
$\ell_{24}$	0.1	0.002	0.006	5.8 %
$\ell_{25}$	0.23	0.001	0.004	1.9 %
$\ell_J$	<b>0.1692</b>	<b>0.0001</b>	<b>0.0042</b>	<b>2.48 %</b>



### 5.3.2. Modelo RDINA con covariable afectando solo ítems

La tabla 5.6 presenta los resultados de la precisión en la recuperación de parámetros del vector  $\mathbf{b}$  y se observa que el sesgo porcentual de los parámetros fue igual o inferior al 4.8 %. Además, el sesgo porcentual de  $\mathbf{b_K}$  que representa al promedio de las estimaciones de los parámetros fue de 0.6 %.

**Tabla 5.6:** Estimación del vector de parámetros  $\mathbf{b}$  del modelo RDINA con covariable afectando solo ítems en la simulación

Parámetro	Valor real	MSE	Sesgo	% Sesgo
$b_1$	0.14	0.018	0.003	2 %
$b_2$	0.59	0.193	0.021	3.5 %
$b_3$	2	0.358	-0.073	3.7 %
$b_4$	0.94	0.513	-0.045	4.8 %
$b_5$	0.47	0.22	0.018	3.8 %
$b_6$	0.9	0.587	0.022	2.4 %
$b_7$	0.67	0.153	0.022	3.3 %
$\mathbf{b_K}$	<b>0.8157</b>	<b>0.041</b>	<b>-0.005</b>	<b>0.6 %</b>

La tabla 5.7 muestra la recuperación de parámetros del vector  $\mathbf{f}$  y se aprecia que el sesgo porcentual de los parámetros fue igual o inferior al 4.8 %. Además, el sesgo porcentual de  $\mathbf{f_J}$  que representa al promedio de las estimaciones de los parámetros fue de 1 %.

**Tabla 5.7:** Estimación del vector de parámetros  $\mathbf{f}$  del modelo RDINA con covariable afectando solo ítems en la simulación

Parámetro	Valor real	MSE	Sesgo	% Sesgo
$f_1$	-1.49	0.266	0.008	0.5 %
$f_2$	-5.25	2.231	0.114	2.2 %
$f_3$	-3.07	0.488	-0.027	0.9 %
$f_4$	-1.74	0.346	-0.048	2.8 %
$f_5$	-4.19	0.746	0.063	1.5 %
$f_6$	-2.3	1.095	0.087	3.8 %
$f_7$	-2.8	0.611	0.051	1.8 %
$f_8$	-3.04	0.456	0.008	0.3 %
$f_9$	-3.13	1.569	-0.072	2.3 %
$f_{10}$	-2.13	0.514	-0.089	4.2 %
$f_{11}$	-1.56	0.223	0.006	0.4 %
$f_{12}$	-0.85	0.215	0.028	3.3 %
$f_{13}$	-2.33	0.445	0.038	1.6 %
$f_{14}$	-2.94	0.443	-0.037	1.3 %
$f_{15}$	-2.33	0.396	0.024	1 %
$f_{16}$	-5.54	1.329	0.037	0.7 %
$f_{17}$	-4.55	0.874	0.062	1.4 %
$f_{18}$	-3.37	0.551	0.09	2.7 %
$f_{19}$	-4.96	1.053	0.134	2.7 %
$f_{20}$	-1.89	0.276	0.091	4.8 %
$f_{21}$	-4.98	0.989	0.112	2.3 %
$f_{22}$	-1.61	0.25	0.019	1.2 %
$f_{23}$	-4.36	0.827	0.073	1.7 %
$f_{24}$	-3.06	1.378	-0.095	3.1 %
$f_{25}$	-3.04	0.583	0.113	3.7 %
$\mathbf{f_J}$	<b>-3.0604</b>	<b>0.306</b>	<b>0.032</b>	<b>1 %</b>

La tabla 5.8 muestra los resultados en la recuperación de parámetros del vector  $\mathbf{d}$  y se aprecia que el sesgo porcentual de los parámetros fue igual o inferior al 5.1 %, excepto los parámetros  $d_{10}$ ,  $d_{11}$  y  $d_7$  que fueron de 9.1 %, 9.1 %, y 6.9 %, respectivamente. Además, el sesgo porcentual de  $\mathbf{d_J}$  que representa al promedio de las estimaciones de los parámetros fue de 1.5 %.

**Tabla 5.8:** Estimación del vector de parámetros  $\mathbf{d}$  del modelo RDINA con covariable afectando solo ítems en la simulación

Parámetro	Valor real	MSE	Sesgo	% Sesgo
$d_1$	1.41	0.219	-0.009	0.6 %
$d_2$	2.09	0.823	-0.063	3 %
$d_3$	1.99	0.381	0.05	2.5 %
$d_4$	1.3	0.3	0.051	3.9 %
$d_5$	2.34	0.439	-0.03	1.3 %
$d_6$	3.78	4.466	0.192	5.1 %
$d_7$	1.48	1.036	0.101	6.9 %
$d_8$	1.63	0.421	0.049	3 %
$d_9$	2.41	1.252	-0.006	0.2 %
$d_{10}$	1.46	0.662	0.134	9.1 %
$d_{11}$	0.47	0.099	0.043	9.1 %
$d_{12}$	0.89	0.207	0.032	3.6 %
$d_{13}$	1.61	0.493	0.05	3.1 %
$d_{14}$	1.26	0.168	0.04	3.2 %
$d_{15}$	1.38	0.202	-0.017	1.2 %
$d_{16}$	2.26	0.465	-0.049	2.2 %
$d_{17}$	1.62	0.236	0.032	2 %
$d_{18}$	1.61	0.23	-0.03	1.9 %
$d_{19}$	1.97	0.347	-0.019	1 %
$d_{20}$	1.68	0.318	0.057	3.4 %
$d_{21}$	1.08	0.102	-0.031	2.9 %
$d_{22}$	0.51	0.219	0	0 %
$d_{23}$	2.13	0.358	-0.051	2.4 %
$d_{24}$	1.71	0.541	0.041	2.4 %
$d_{25}$	2.41	0.706	0.075	3.1 %
$\mathbf{d_J}$	<b>1.6992</b>	<b>0.139</b>	<b>0.026</b>	<b>1.5 %</b>

La tabla 5.9 presenta la recuperación de parámetros del vector  $\ell$  y se muestra que el sesgo porcentual de los parámetros fue igual o inferior al 4.7 %. Además, el sesgo porcentual de  $\ell_J$  que representa al promedio de las estimaciones de los parámetros fue de 0.1 %.

**Tabla 5.9:** Estimación del vector de parámetros  $\ell$  del modelo RDINA con covariable afectando solo ítems en la simulación

Parámetro	Valor real	MSE	Sesgo	% Sesgo
$\ell_1$	0.15	0.001	0.001	0.6 %
$\ell_2$	0.19	0.002	-0.002	0.8 %
$\ell_3$	0.12	0.001	0.003	2.5 %
$\ell_4$	0.13	0.001	0.004	2.7 %
$\ell_5$	0.23	0.002	0	0.2 %
$\ell_6$	0.27	0.004	0.001	0.5 %
$\ell_7$	0.15	0.001	0.001	0.6 %
$\ell_8$	0.14	0.001	0.001	0.9 %
$\ell_9$	0.19	0.002	0.005	2.6 %
$\ell_{10}$	0.1	0.001	0.005	4.7 %
$\ell_{11}$	0.05	0	0	0.3 %
$\ell_{12}$	0.11	0.001	-0.002	2 %
$\ell_{13}$	0.18	0.001	-0.001	0.4 %
$\ell_{14}$	0.13	0.001	0.003	2.3 %
$\ell_{15}$	0.17	0.001	-0.001	0.3 %
$\ell_{16}$	0.27	0.002	0	0.1 %
$\ell_{17}$	0.15	0.001	-0.002	1.3 %
$\ell_{18}$	0.19	0.001	-0.003	1.5 %
$\ell_{19}$	0.23	0.002	-0.005	2.3 %
$\ell_{20}$	0.11	0.001	-0.005	4.4 %
$\ell_{21}$	0.16	0.001	-0.002	1.3 %
$\ell_{22}$	0.07	0	0	0.4 %
$\ell_{23}$	0.23	0.002	-0.002	0.8 %
$\ell_{24}$	0.09	0.001	0.001	1.1 %
$\ell_{25}$	0.22	0.002	-0.006	2.6 %
$\ell_J$	<b>0.1612</b>	<b>0</b>	<b>0</b>	<b>0.1 %</b>

### 5.3.3. Modelo RDINA con covariable afectando solo atributos

La tabla 5.10 presenta los resultados de la precisión en la recuperación de parámetros del vector  $\mathbf{b}$  y muestra que el sesgo porcentual de los parámetros fue igual o inferior al 5.6 %, excepto el parámetro  $b_4$  que fue de 8.7 %. Además, el sesgo porcentual de  $\mathbf{b_K}$  que representa al promedio de las estimaciones de los parámetros fue de 3.3 %.

**Tabla 5.10:** Estimación del vector de parámetros  $\mathbf{b}$  del modelo RDINA con covariable afectando solo atributos en la simulación

Parámetro	Valor real	MSE	Sesgo	% Sesgo
$b_1$	-5.67	0.49	-0.173	3 %
$b_2$	-5.75	1.36	-0.174	3 %
$b_3$	-2.6	1.684	-0.033	1.3 %
$b_4$	-1.5	4.353	-0.131	8.7 %
$b_5$	-3.7	0.807	-0.097	2.6 %
$b_6$	-6.33	7.392	-0.352	5.6 %
$b_7$	-3.36	1.661	-0.001	0 %
$\mathbf{b_K}$	<b>-4.13</b>	<b>0.2722</b>	<b>-0.1372</b>	<b>3.3 %</b>

La tabla 5.11 presenta la recuperación de parámetros del vector  $\mathbf{h}$  y muestra que el sesgo porcentual de los parámetros fue igual o inferior al 5.3 %. Además, el sesgo porcentual de  $\mathbf{h_K}$  que representa al promedio de las estimaciones de los parámetros fue de 3.6 %.

La tabla 5.12 muestra la recuperación de parámetros del vector  $\mathbf{f}$  y se aprecia que el sesgo porcentual de los parámetros fue igual o inferior al 4.8 %. Además, el sesgo porcentual de  $\mathbf{f_J}$  que representa al promedio de las estimaciones de los parámetros fue de 1.5 %.

**Tabla 5.11:** Estimación del vector de parámetros  $\mathbf{h}$  del modelo RDINA con covariable afectando solo atributos en la simulación

Parámetro	Valor real	MSE	Sesgo	% Sesgo
$h_1$	0.34	0.002	0.01	2.8 %
$h_2$	0.34	0.004	0.008	2.4 %
$h_3$	0.21	0.006	0.008	3.7 %
$h_4$	0.14	0.009	0.011	8 %
$h_5$	0.24	0.003	0.006	2.4 %
$h_6$	0.4	0.022	0.021	5.3 %
$h_7$	0.22	0.005	0.004	1.7 %
$\mathbf{h_K}$	<b>0.27</b>	<b>0.0008</b>	<b>0.0096</b>	<b>3.6 %</b>

**Tabla 5.12:** Estimación del vector de parámetros  $\mathbf{f}$  del modelo RDINA con covariable afectando solo atributos en la simulación

Parámetro	Valor real	MSE	Sesgo	% Sesgo
$f_1$	0.89	0.023	-0.025	2.8 %
$f_2$	-2.24	0.316	-0.107	4.8 %
$f_3$	-0.77	0.015	-0.005	0.6 %
$f_4$	0.21	0.013	-0.013	6 %
$f_5$	-0.45	0.019	0.017	3.8 %
$f_6$	1.53	0.034	0.031	2 %
$f_7$	-0.57	0.019	0.013	2.3 %
$f_8$	-0.92	0.015	0.004	0.4 %
$f_9$	-0.27	0.044	-0.001	0.2 %
$f_{10}$	-0.69	0.033	0.023	3.3 %
$f_{11}$	-0.88	0.024	-0.001	0.1 %
$f_{12}$	0.84	0.016	-0.013	1.5 %
$f_{13}$	0.52	0.016	0.01	1.9 %
$f_{14}$	-1.02	0.019	-0.019	1.9 %
$f_{15}$	0.21	0.022	0.002	0.8 %
$f_{16}$	-1.5	0.042	0.007	0.5 %
$f_{17}$	-1.94	0.038	0.013	0.7 %
$f_{18}$	-0.41	0.019	0.007	1.6 %
$f_{19}$	-1.37	0.028	-0.012	0.8 %
$f_{20}$	-0.11	0.015	0.003	2.3 %
$f_{21}$	-2.7	0.087	-0.085	3.1 %
$f_{22}$	-0.66	0.025	0.002	0.3 %
$f_{23}$	-0.8	0.032	-0.007	0.9 %
$f_{24}$	-1.56	0.117	-0.043	2.8 %
$f_{25}$	0.44	0.018	-0.008	1.8 %
$\mathbf{f_J}$	<b>-0.5688</b>	<b>0.0024</b>	<b>-0.0083</b>	<b>1.5 %</b>

La tabla 5.13 muestra los resultados en la recuperación de parámetros del vector  $\mathbf{d}$  y se aprecia que el sesgo porcentual de los parámetros fue igual o inferior al 3.9 %. Además, el sesgo porcentual de  $\mathbf{d_J}$  que representa al promedio de las estimaciones de los parámetros fue de 1.7 %.

**Tabla 5.13:** Estimación del vector de parámetros  $\mathbf{d}$  del modelo RDINA con covariable afectando solo atributos en la simulación

Parámetro	Valor real	MSE	Sesgo	% Sesgo
$d_1$	1.48	0.074	0.031	2.1 %
$d_2$	3.12	0.283	0.098	3.1 %
$d_3$	1.62	0.055	0.009	0.5 %
$d_4$	2.06	0.142	0.081	3.9 %
$d_5$	2.83	0.132	0.06	2.1 %
$d_6$	5.27	2.658	0.037	0.7 %
$d_7$	2.23	0.397	0.065	2.9 %
$d_8$	1.93	0.055	-0.012	0.6 %
$d_9$	2.83	0.189	0.071	2.5 %
$d_{10}$	1.72	0.121	0.018	1 %
$d_{11}$	0.56	0.061	0.017	3.1 %
$d_{12}$	1.34	0.077	0.039	2.9 %
$d_{13}$	2.08	0.193	0.071	3.4 %
$d_{14}$	2.01	0.074	0.032	1.6 %
$d_{15}$	2.25	0.086	0.023	1 %
$d_{16}$	3.04	0.081	0.012	0.4 %
$d_{17}$	1.7	0.061	0.022	1.3 %
$d_{18}$	2.37	0.104	0.061	2.6 %
$d_{19}$	2.77	0.071	-0.015	0.5 %
$d_{20}$	2.24	0.155	0.014	0.6 %
$d_{21}$	1.94	0.084	0.075	3.9 %
$d_{22}$	1.11	0.068	0.025	2.2 %
$d_{23}$	2.75	0.078	-0.008	0.3 %
$d_{24}$	1.86	0.146	0.06	3.2 %
$d_{25}$	2.75	0.301	0.046	1.7 %
$\mathbf{d_J}$	<b>2.2344</b>	<b>0.0111</b>	<b>0.0372</b>	<b>1.7 %</b>

#### 5.3.4. Modelo RDINA sin covariable

La tabla 5.14 presenta los resultados de la precisión en la recuperación de parámetros del vector  $\mathbf{b}$  y muestra que el sesgo porcentual de los parámetros fue igual o inferior al 5.7 %. Además, el sesgo porcentual de  $\mathbf{b_K}$  que representa al promedio de las estimaciones de los parámetros fue de 0.9 %.

**Tabla 5.14:** Estimación del vector de parámetros  $\mathbf{b}$  del modelo RDINA sin covariable en la simulación

Parámetro	Valor real	MSE	Sesgo	% Sesgo
$b_1$	0.1	0.011	0.001	1.4 %
$b_2$	0.63	0.055	-0.036	5.7 %
$b_3$	1.69	0.124	-0.066	3.9 %
$b_4$	1.51	0.494	-0.01	0.7 %
$b_5$	0.45	0.075	0.022	4.9 %
$b_6$	0.95	0.236	0.015	1.6 %
$b_7$	0.75	0.051	0.017	2.2 %
$\mathbf{b_K}$	<b>0.8686</b>	<b>0.0255</b>	<b>-0.0081</b>	<b>0.9 %</b>

La tabla 5.15 muestra la recuperación de parámetros del vector  $\mathbf{f}$  y se aprecia que el sesgo porcentual de los parámetros fue igual o inferior al 4.8 %, excepto el parámetro  $f_9$  que fue de 6.4 %. Además, el sesgo porcentual de  $\mathbf{f_J}$  que representa al promedio de las estimaciones de los parámetros fue de 0.3 %.

La tabla 5.16 muestra los resultados en la recuperación de parámetros del vector  $\mathbf{d}$  y se aprecia que el sesgo porcentual de los parámetros fue igual o inferior al 5.8 %, excepto el parámetro  $d_7$  que fue de 7.6 %. Además, el sesgo porcentual de  $\mathbf{d_J}$  que representa al promedio



**Tabla 5.15:** Estimación del vector de parámetros  $\mathbf{f}$  del modelo RDINA sin covariable en la simulación

Parámetro	Valor real	MSE	Sesgo	% Sesgo
$f_1$	0.91	0.022	-0.004	0.5 %
$f_2$	-2.43	0.603	0.057	2.4 %
$f_3$	-0.94	0.018	0.007	0.8 %
$f_4$	0.14	0.018	-0.007	4.8 %
$f_5$	-0.42	0.014	0.017	4.2 %
$f_6$	1.31	0.046	0.016	1.2 %
$f_7$	-0.71	0.038	-0.03	4.2 %
$f_8$	-0.88	0.015	0.019	2.2 %
$f_9$	-0.24	0.051	-0.015	6.4 %
$f_{10}$	-0.79	0.05	-0.01	1.2 %
$f_{11}$	-0.88	0.024	-0.009	1 %
$f_{12}$	0.82	0.022	-0.002	0.2 %
$f_{13}$	0.49	0.015	0.012	2.5 %
$f_{14}$	-1.08	0.019	0.01	0.9 %
$f_{15}$	0.29	0.023	-0.007	2.2 %
$f_{16}$	-1.37	0.038	-0.011	0.8 %
$f_{17}$	1.96	0.029	0.023	1.2 %
$f_{18}$	-0.41	0.014	0.012	2.9 %
$f_{19}$	-1.32	0.024	-0.018	1.3 %
$f_{20}$	-0.11	0.012	0.003	2.8 %
$f_{21}$	-2.58	0.083	-0.038	1.5 %
$f_{22}$	-0.61	0.027	-0.019	3.2 %
$f_{23}$	-0.74	0.023	0	0 %
$f_{24}$	-1.62	0.188	-0.044	2.7 %
$f_{25}$	0.42	0.012	-0.001	0.1 %
$\mathbf{f_J}$	<b>-0.5884</b>	<b>0.0031</b>	<b>-0.0015</b>	<b>0.3 %</b>

de las estimaciones de los parámetros fue de 1.2 %.

**Tabla 5.16:** Estimación del vector de parámetros  $\mathbf{d}$  del modelo RDINA sin covariable en la simulación

Parámetro	Valor real	MSE	Sesgo	% Sesgo
$d_1$	1.51	0.076	0.015	1 %
$d_2$	2.58	0.792	-0.061	2.4 %
$d_3$	1.78	0.093	0.028	1.5 %
$d_4$	1.93	0.190	0.089	4.6 %
$d_5$	2.75	0.134	0.009	0.3 %
$d_6$	4.6	2.022	-0.266	5.8 %
$d_7$	1.87	0.273	0.143	7.6 %
$d_8$	1.99	0.116	0.029	1.5 %
$d_9$	2.62	0.236	0.102	3.9 %
$d_{10}$	1.65	0.1	0.053	3.2 %
$d_{11}$	0.53	0.052	0.029	5.4 %
$d_{12}$	1.28	0.132	0.032	2.5 %
$d_{13}$	2.14	0.147	0.046	2.2 %
$d_{14}$	1.91	0.127	0.057	3 %
$d_{15}$	2.13	0.116	0.019	0.9 %
$d_{16}$	2.98	0.094	0.03	1 %
$d_{17}$	1.74	0.057	-0.018	1 %
$d_{18}$	2.35	0.114	0.056	2.4 %
$d_{19}$	2.65	0.109	0.07	2.6 %
$d_{20}$	2.07	0.141	-0.016	0.8 %
$d_{21}$	1.81	0.092	0.039	2.2 %
$d_{22}$	0.91	0.074	0.023	2.6 %
$d_{23}$	2.78	0.059	0.017	0.6 %
$d_{24}$	1.91	0.238	0.047	2.5 %
$d_{25}$	2.49	0.21	0.045	1.8 %
$\mathbf{d_J}$	<b>2.1184</b>	<b>0.0097</b>	<b>0.0247</b>	<b>1.2 %</b>

## Capítulo 6

# Aplicación

La aplicación de este trabajo se orienta al ámbito de la educación superior, específicamente a una prueba de admisión rendida en una universidad privada del Perú por 727 jóvenes que deseaban estudiar carreras orientada a las especialidades de Letras y Arquitectura en el año 2018.

En el presente capítulo se examina el ajuste de la extensión del modelo RDINA usando una covariable con la finalidad de analizar su efecto en los siguientes modelos:

- a) Modelo RDINA con covariable afectando a los ítems y atributos
- b) Modelo RDINA con covariable afectando solo a los ítems
- c) Modelo RDINA con covariable afectando solo a los atributos
- d) Modelo RDINA sin covariable.

Cabe indicar que dicha prueba de admisión evalúa las siguientes competencias consideradas de gran importancia para lograr un buen desempeño durante la etapa universitaria: Lectura, Redacción y Matemática, las cuales debieron desarrollarse durante toda la etapa escolar. La Universidad espera que los resultados en las pruebas de admisión permitan predecir el desempeño académico de los nuevos alumnos admitidos mediante esta prueba, principalmente, en el primer año de estudios universitarios y disminuir la deserción originadas por el bajo rendimiento. La prueba de admisión consta de 120 preguntas de opción múltiple y tiene una duración de 3 horas y 15 minutos.

A continuación se muestran algunos detalles de las evaluaciones de las competencias en Lectura, Redacción y Matemática:

**Tabla 6.1:** Competencias evaluadas en la prueba de admisión en la aplicación

Sección	Competencias	N° Preguntas	Tiempo	Peso
1	Lectura	36	1 hora	25 %
2	Redacción	36	40 minutos	25 %
3	Matemática	48	1 hora, 35 minutos	50 %

Asimismo con el apoyo de una profesional en Ciencias Lingüísticas se elaboró una matriz  $Q$  de habilidad para la competencia Redacción que será objeto de estudio. Estas tienen una dimensión  $J \times K$ , siendo  $J$  la cantidad total de ítems ( $J = 36$ ) y  $K$  es la cantidad de habilidades ( $K = 8$ ).

Las habilidades definidas por la profesional para la competencia Redacción son las siguientes:

**Tabla 6.2:** Habilidades definidas en la prueba de admisión en la competencia Redacción

Habilidad	Descripción de las habilidades
$R_1$	Interpreta término / expresión a partir del co-texto
$R_2$	Identifica significado de palabra
$R_3$	Conoce reglas morfo-sintácticas
$R_4$	Conoce reglas ortográficas
$R_5$	Conoce reglas de puntuación
$R_6$	Reconoce organización textual
$R_7$	Comprende texto de manera global
$R_8$	Usa / identifica habilidades discursivas académicas

Asimismo la universidad que aplico la prueba de admisión en la competencia Redacción lo dividió en los siguientes temas y subtemas:

**Tabla 6.3:** División de la prueba de admisión en la competencia Redacción por temas y subtemas

Temas y subtemas en Redacción	Preguntas
<b>Ortografía y puntuación</b>	
Utilización correcta de la grafía	4, 8, 17, 18, 25, 29, 34
Uso correcto de la tilde	10
Uso correcto de las reglas de puntuación	19, 20, 32, 33
<b>Vocabulario y construcción oracional</b>	
Reconoce la pertinencia del uso del vocabulario	1, 3, 15, 21
Determina la pertinencia del uso de conectores	36
Reconoce si un texto propuesto es gramaticalmente correcto	2, 22, 23, 26, 30, 31, 35
<b>Contenido y organización lógica de ideas</b>	
Reconoce la cohesión temática de textos	7, 11, 12, 14, 24, 27
Reconoce el orden lógico de las ideas	5, 6, 9, 13, 16, 28

Se usarán las respuestas de los 727 jóvenes que rindieron la prueba de admisión como matriz de respuestas dicotómicas.

El objetivo de este estudio será evaluar el ajuste de los modelos bajo los 4 escenarios indicados, estimar sus parámetros mediante los algoritmos EM y Newton-Raphson y analizar el efecto de la covariable en cada uno de ellos, para lo cual se empleará el paquete Latent Gold.

**Tabla 6.4:** Número de parámetros a estimar según modelo

Modelo	N° de parámetros
RDINA con covariable afectando ítems y atributos	124
RDINA con covariable afectando solo ítems	116
RDINA con covariable afectando solo atributos	88
RDINA sin covariable	80

La matriz **Q** de la competencia **Redacción** es la siguiente:

**Tabla 6.5:** Matriz Q de habilidades de la competencia Redacción

Ítem	R <sub>1</sub>	R <sub>2</sub>	R <sub>3</sub>	R <sub>4</sub>	R <sub>5</sub>	R <sub>6</sub>	R <sub>7</sub>	R <sub>8</sub>	Total
1	1	1	0	0	0	0	0	0	2
2	0	0	1	0	0	0	0	0	1
3	1	1	1	0	0	0	0	0	3
4	0	0	0	1	0	0	0	0	1
5	0	0	0	0	0	1	1	1	3
6	0	0	0	0	0	1	1	1	3
7	0	0	0	0	0	1	1	1	3
8	0	0	0	1	0	0	0	0	1
9	0	0	0	0	0	1	1	1	3
10	1	0	0	1	0	0	0	0	2
11	0	0	0	0	0	1	1	1	3
12	0	0	0	0	0	1	1	1	3
13	0	0	0	0	0	1	1	1	3
14	0	0	0	0	0	1	1	1	3
15	1	1	0	0	0	0	1	0	3
16	0	0	0	0	0	1	1	1	3
17	1	0	0	1	0	0	0	0	2
18	1	0	0	1	0	0	0	0	2
19	0	0	0	0	1	0	1	0	2
20	0	0	0	0	1	0	1	0	2
21	1	1	0	0	0	0	0	0	2
22	0	0	1	0	0	0	0	0	1
23	0	0	1	0	0	0	0	0	1
24	0	0	0	0	0	1	1	1	3
25	0	0	0	1	0	0	0	0	1
26	0	0	1	0	0	0	0	0	1
27	0	0	0	0	0	1	1	1	3
28	0	0	0	0	0	1	1	1	3
29	1	0	0	1	0	0	0	0	2
30	0	0	1	0	0	0	0	0	1
31	0	0	1	0	0	0	0	0	1
32	0	0	0	0	1	0	1	1	3
33	0	0	0	0	1	0	1	1	3
34	1	0	0	1	0	0	0	0	2
35	0	0	1	0	0	0	0	0	1
36	0	0	0	0	0	1	1	1	3
Total	9	4	8	8	4	13	18	15	79

Para la aplicación del modelo RDINA con covariable, además de tener la matriz *Q* se cuenta con las respuestas dicotómicas de los 727 examinados de cada una de los 36 ítems en la competencia Redacción, así como también el puntaje obtenido en la competencia Lectura.

En el apéndice C se incluyen los códigos para la estimación de los 4 modelos indicados.

A continuación se muestran resultados del ajuste de los modelos indicados así como la estimación de sus parámetros teniendo en cuenta lo siguiente:

- La **matriz de respuestas dicotómicas** que se estudió corresponde a las respuestas de todos los examinados para cada uno de los ítems en la prueba de admisión en la

competencia **Redacción**.

- La **covariable** que se usó fue el puntaje general obtenido por cada uno de los examinados en la misma prueba de admisión en la competencia **Lectura**.

En la tabla 6.6 se muestra los resultados de ajuste, de los cuatro modelos utilizados para este estudio. Para seleccionar el modelo que mejor se ajuste, se utilizaron los criterios AIC y BIC. Si consideramos el AIC, el modelo RDINA con covariable afectando ítems y atributos se ajustó mejor al obtener el índice más bajo, mientras si consideramos el BIC, el modelo RDINA con covariable afectando solo atributos se ajustó mejor, dado que su índice fue el más bajo. Por otro lado, para ambos criterios el modelo RDINA sin covariable resultó el modelo con el menor ajuste, y el modelo RDINA con covariable afectando solo ítems resultó el segundo con mejor ajuste.

**Tabla 6.6:** Resultados generales de los modelos ajustados en la aplicación

Modelo	N° parámetros	BIC (LL)	AIC (LL)	Log verosimilitud	Log priori	Log posterior
RDINA con covariable afectando ítems y atributos	<b>124</b>	27,409.88	26,840.85	-13,296.42	-31.15	-13,327.58
RDINA con covariable afectando solo ítems	<b>116</b>	27,388.95	26,856.64	-13,312.32	-30.57	-13,342.89
RDINA con covariable afectando solo atributos	<b>88</b>	27,342.58	26,938.75	-13,381.38	-31.77	-13,413.15
RDINA sin covariable	<b>80</b>	27,598.26	27,231.14	-13,535.57	-29.72	-13,565.29

## 6.1. Estimación de parámetros del modelo RDINA con covariable afectando ítems y atributos

La tabla 6.7 presenta las estimaciones de los parámetros a nivel de atributos de los vectores **b** y **h**. En el caso del vector **b**, las estimaciones de los parámetros  $b_3$  y  $b_4$  fueron significativas con un valor  $p < 0.05$ , lo cual quiere decir, que los examinados con bajos niveles de Lectura tienden a tener un mayor dominio de las habilidades **R<sub>3</sub>** y **R<sub>4</sub>**. En el vector **h**, las estimaciones de los parámetros  $h_3$ ,  $h_4$  y  $h_5$  tuvieron un valor  $p < 0.05$  con valores positivos, por tanto, el puntaje en Lectura tuvo un efecto positivo en la probabilidad de que un examinado posea la habilidades **R<sub>3</sub>**, **R<sub>4</sub>** y **R<sub>5</sub>**, que están referidas al conocimiento de las reglas morfo-sintácticas, las reglas ortográficas y las reglas de puntuación, respectivamente.

**Tabla 6.7:** Estimación de los vectores de parámetros **b** y **h** del modelo RDINA con covariable afectando ítems y atributos en la aplicación

Vector de parámetros b				Vector de parámetros h			
Param.	Est.	s.e.	p-valor	Param.	Est.	s.e.	p-valor
$b_1$	-0.186	2.153	0.93	$h_1$	0.003	0.003	0.44
$b_2$	-1.779	2.358	0.45	$h_2$	0.004	0.004	0.34
$b_3$	-4.837	2.38	0.042	$h_3$	0.008	0.004	0.044
$b_4$	-5.48	1.631	0.001	$h_4$	0.01	0.003	0
$b_5$	-7.208	4.522	0.11	$h_5$	0.018	0.009	0.048
$b_6$	-6.638	4.756	0.16	$h_6$	0.017	0.009	0.068
$b_7$	-0.122	1.763	0.94	$h_7$	0.002	0.003	0.58
$b_8$	-0.092	5.271	0.99	$h_8$	0.005	0.009	0.62

La tabla 6.8 presenta las estimaciones de los parámetros a nivel de ítems de los vectores **f**, **d** y **ℓ**. Los parámetros que tuvieron un valor  $p < 0.05$  fueron 17 en el caso del vector **f**, con



valores negativos, en lo que respecta al vector  $\mathbf{d}$ , fueron 28 parámetros con valores positivos, y en el caso del vector  $\ell$ , fueron 21 parámetros con valores positivos.

**Tabla 6.8:** Estimación de los vectores de parámetros  $\mathbf{f}$ ,  $\mathbf{d}$  y  $\ell$  del modelo RDINA con covariable afectando ítems y atributos en la aplicación

Vector de parámetros $\mathbf{f}$				Vector de parámetros $\mathbf{d}$				Vector de parámetros $\ell$			
Param.	Est.	s.e.	p-valor	Param.	Est.	s.e.	p-valor	Param.	Est.	s.e.	p-valor
$f_1$	-1.052	0.895	0.24	$d_1$	1.072	0.454	0.018	$\ell_1$	0.003	0.002	0.03
$f_2$	-4.837	2.253	0.032	$d_2$	1.677	2.587	0.52	$\ell_2$	0.015	0.004	0.001
$f_3$	0.612	1.053	0.56	$d_3$	4.555	3.945	0.25	$\ell_3$	-0.001	0.002	0.49
$f_4$	-2.649	0.717	0	$d_4$	0.75	0.262	0.004	$\ell_4$	0.005	0.001	0.001
$f_5$	-0.812	0.661	0.22	$d_5$	0.603	0.243	0.013	$\ell_5$	0.002	0.001	0.095
$f_6$	-0.955	0.627	0.13	$d_6$	0.474	0.243	0.051	$\ell_6$	0.002	0.001	0.091
$f_7$	-0.271	1.101	0.81	$d_7$	1.673	0.439	0	$\ell_7$	0.003	0.002	0.16
$f_8$	-3.797	1.019	0	$d_8$	1.022	0.337	0.002	$\ell_8$	0.009	0.002	0
$f_9$	-1.392	0.966	0.15	$d_9$	0.981	0.324	0.003	$\ell_9$	0.005	0.002	0.006
$f_{10}$	-2.608	0.924	0.005	$d_{10}$	1.686	0.343	0	$\ell_{10}$	0.006	0.002	0.001
$f_{11}$	-3.422	0.732	0	$d_{11}$	0.741	0.242	0.002	$\ell_{11}$	0.006	0.001	0
$f_{12}$	-1.069	1.024	0.3	$d_{12}$	1.358	0.397	0.001	$\ell_{12}$	0.004	0.002	0.028
$f_{13}$	-0.477	0.6	0.43	$d_{13}$	0.015	0.225	0.95	$\ell_{13}$	0	0.001	0.95
$f_{14}$	-0.319	1.234	0.8	$d_{14}$	2.335	0.567	0	$\ell_{14}$	0.003	0.002	0.19
$f_{15}$	-1.947	0.713	0.006	$d_{15}$	0.937	0.396	0.018	$\ell_{15}$	0.004	0.001	0.004
$f_{16}$	-1.669	1.201	0.16	$d_{16}$	1.245	0.427	0.004	$\ell_{16}$	0.006	0.002	0.005
$f_{17}$	-2.665	1.159	0.022	$d_{17}$	3.264	0.764	0	$\ell_{17}$	0.004	0.002	0.043
$f_{18}$	-3.979	0.959	0	$d_{18}$	1.418	0.332	0	$\ell_{18}$	0.009	0.002	0
$f_{19}$	-3.252	1.314	0.013	$d_{19}$	2.064	0.432	0	$\ell_{19}$	0.006	0.002	0.005
$f_{20}$	-1.064	1.425	0.46	$d_{20}$	2.272	0.705	0.001	$\ell_{20}$	0.004	0.003	0.092
$f_{21}$	-3.524	1.232	0.004	$d_{21}$	1.919	0.704	0.006	$\ell_{21}$	0.006	0.002	0.004
$f_{22}$	-0.626	1.408	0.66	$d_{22}$	2.236	0.843	0.008	$\ell_{22}$	0.001	0.003	0.75
$f_{23}$	-0.01	0.801	0.99	$d_{23}$	0.936	0.324	0.004	$\ell_{23}$	-0.002	0.001	0.089
$f_{24}$	-3.419	0.88	0	$d_{24}$	1.818	0.335	0	$\ell_{24}$	0.003	0.001	0.079
$f_{25}$	-1.696	1.013	0.094	$d_{25}$	1.83	0.377	0	$\ell_{25}$	0.003	0.002	0.072
$f_{26}$	0.026	0.757	0.97	$d_{26}$	-0.008	0.34	0.98	$\ell_{26}$	0.002	0.001	0.12
$f_{27}$	-3.148	1.298	0.015	$d_{27}$	2.46	0.772	0.001	$\ell_{27}$	0.007	0.002	0.001
$f_{28}$	-4.414	0.888	0	$d_{28}$	1.168	0.278	0	$\ell_{28}$	0.008	0.002	0
$f_{29}$	-2.539	0.95	0.008	$d_{29}$	2.334	0.331	0	$\ell_{29}$	0.003	0.002	0.042
$f_{30}$	-1.887	0.889	0.034	$d_{30}$	1.265	0.331	0	$\ell_{30}$	0.001	0.002	0.58
$f_{31}$	-1.014	0.805	0.21	$d_{31}$	-0.145	0.339	0.67	$\ell_{31}$	0.004	0.002	0.004
$f_{32}$	-3.153	0.955	0.001	$d_{32}$	0.436	0.325	0.18	$\ell_{32}$	0.008	0.002	0
$f_{33}$	-0.289	0.84	0.73	$d_{33}$	1.273	0.283	0	$\ell_{33}$	0.001	0.001	0.56
$f_{34}$	-1.776	1.08	0.1	$d_{34}$	1.976	0.456	0	$\ell_{34}$	0.005	0.002	0.007
$f_{35}$	-2.607	0.884	0.003	$d_{35}$	1.072	0.403	0.008	$\ell_{35}$	0.005	0.002	0
$f_{36}$	-0.288	0.595	0.63	$d_{36}$	0.392	0.227	0.084	$\ell_{36}$	0	0.001	0.83

## 6.2. Estimación de parámetros del modelo RDINA con covariable afectando solo ítems

La tabla 6.9 presenta las estimaciones de los parámetros a nivel de atributos del vector  $\mathbf{b}$ . Como puede apreciarse, las estimaciones de los parámetros  $b_1$  y  $b_6$  fueron positivos y significativas con un valor  $p < 0.05$ , lo cual quiere decir, que los examinados tienden a tener un mayor dominio de las habilidades  $\mathbf{R}_1$  y  $\mathbf{R}_6$ , referidas a la interpretación del término/expresión a partir del co-texto y al reconocimiento de la organización textual, respectivamente.

La tabla 6.10 presenta las estimaciones de los parámetros a nivel de ítems de los vectores  $\mathbf{f}$ ,  $\mathbf{d}$  y  $\ell$ . Los parámetros que tuvieron un valor  $p < 0.05$  fueron 19 en el caso del vector  $\mathbf{f}$ , con valores negativos, en lo que respecta al vector  $\mathbf{d}$ , fueron 26 parámetros con valores positivos y negativos, y en el caso del vector  $\ell$ , fueron 30 parámetros con valores positivos.

**Tabla 6.9:** Estimación del vector de parámetros **b** del modelo RDINA con covariable afectando solo ítems en la aplicación

Param.	Est.	s.e.	p-valor
$b_1$	1.549	0.535	0.004
$b_2$	0.218	0.466	0.64
$b_3$	-0.58	0.55	0.29
$b_4$	0.115	0.348	0.74
$b_5$	2.641	1.359	0.052
$b_6$	2.566	1.243	0.039
$b_7$	-0.336	0.318	0.29
$b_8$	2.335	1.24	0.06

**Tabla 6.10:** Estimación de los vectores de parámetros **f**, **d** y **l** del modelo RDINA con covariable afectando solo ítems en la aplicación

Vector de parámetros f				Vector de parámetros d				Vector de parámetros l			
Param.	Est.	s.e.	p-valor	Param.	Est.	s.e.	p-valor	Param.	Est.	s.e.	p-valor
$f_1$	-1.384	0.795	0.082	$d_1$	0.668	0.356	0.061	$\ell_1$	0.004	0.001	0.002
$f_2$	-5.911	2.203	0.007	$d_2$	1.298	1.447	0.37	$\ell_2$	0.017	0.004	0
$f_3$	-1.232	0.688	0.073	$d_3$	3.09	2.484	0.21	$\ell_3$	0.002	0.001	0.072
$f_4$	-3.38	0.694	0	$d_4$	0.688	0.262	0.009	$\ell_4$	0.006	0.001	0
$f_5$	-0.528	0.651	0.42	$d_5$	-0.601	0.252	0.017	$\ell_5$	0.002	0.001	0.028
$f_6$	-0.706	0.624	0.26	$d_6$	-0.522	0.25	0.037	$\ell_6$	0.002	0.001	0.035
$f_7$	0.46	1.061	0.66	$d_7$	-1.63	0.405	0	$\ell_7$	0.004	0.002	0.019
$f_8$	-4.761	0.974	0	$d_8$	0.801	0.334	0.016	$\ell_8$	0.011	0.002	0
$f_9$	-0.935	0.952	0.33	$d_9$	-1.058	0.321	0.001	$\ell_9$	0.006	0.002	0.001
$f_{10}$	-4.016	0.885	0	$d_{10}$	1.609	0.36	0	$\ell_{10}$	0.008	0.002	0
$f_{11}$	-3.082	0.707	0	$d_{11}$	-0.811	0.25	0.001	$\ell_{11}$	0.007	0.001	0
$f_{12}$	-0.538	0.994	0.59	$d_{12}$	-1.254	0.386	0.001	$\ell_{12}$	0.005	0.002	0.002
$f_{13}$	-0.449	0.609	0.46	$d_{13}$	-0.06	0.233	0.8	$\ell_{13}$	0	0.001	0.95
$f_{14}$	0.656	1.17	0.57	$d_{14}$	-2.121	0.451	0	$\ell_{14}$	0.005	0.002	0.016
$f_{15}$	-1.912	0.643	0.003	$d_{15}$	-0.425	0.365	0.24	$\ell_{15}$	0.004	0.001	0
$f_{16}$	-1.15	1.196	0.34	$d_{16}$	-1.222	0.417	0.003	$\ell_{16}$	0.007	0.002	0.001
$f_{17}$	-5.411	1.001	0	$d_{17}$	2.93	0.527	0	$\ell_{17}$	0.009	0.002	0
$f_{18}$	-5.182	0.952	0	$d_{18}$	1.411	0.36	0	$\ell_{18}$	0.011	0.002	0
$f_{19}$	-2.366	0.964	0.014	$d_{19}$	-1.893	0.402	0	$\ell_{19}$	0.008	0.002	0
$f_{20}$	-0.007	1.291	1	$d_{20}$	-2.112	0.614	0.001	$\ell_{20}$	0.006	0.002	0.003
$f_{21}$	-7.055	2.84	0.013	$d_{21}$	4.359	2.982	0.14	$\ell_{21}$	0.011	0.004	0.008
$f_{22}$	-2.536	0.798	0.002	$d_{22}$	1.825	0.628	0.004	$\ell_{22}$	0.005	0.001	0.001
$f_{23}$	-0.919	0.721	0.2	$d_{23}$	1.096	0.334	0.001	$\ell_{23}$	-0.001	0.001	0.46
$f_{24}$	-2.479	0.72	0.001	$d_{24}$	-1.841	0.373	0	$\ell_{24}$	0.004	0.001	0.001
$f_{25}$	-3.661	0.889	0	$d_{25}$	1.652	0.377	0	$\ell_{25}$	0.007	0.002	0
$f_{26}$	0.056	0.74	0.94	$d_{26}$	-0.079	0.364	0.83	$\ell_{26}$	0.002	0.001	0.084
$f_{27}$	-2.299	1.197	0.055	$d_{27}$	-2.186	0.547	0	$\ell_{27}$	0.01	0.002	0
$f_{28}$	-3.953	0.813	0	$d_{28}$	-1.074	0.274	0	$\ell_{28}$	0.009	0.001	0
$f_{29}$	-4.91	0.895	0	$d_{29}$	2.372	0.396	0	$\ell_{29}$	0.008	0.002	0
$f_{30}$	-3.187	0.766	0	$d_{30}$	1.406	0.36	0	$\ell_{30}$	0.003	0.001	0.009
$f_{31}$	-0.921	0.782	0.24	$d_{31}$	-0.023	0.345	0.95	$\ell_{31}$	0.004	0.001	0.003
$f_{32}$	-2.993	0.954	0.002	$d_{32}$	-0.452	0.318	0.16	$\ell_{32}$	0.009	0.002	0
$f_{33}$	0.193	0.716	0.79	$d_{33}$	-1.221	0.274	0	$\ell_{33}$	0.002	0.001	0.089
$f_{34}$	-3.290	1.05	0.002	$d_{34}$	2.022	0.5	0	$\ell_{34}$	0.008	0.002	0
$f_{35}$	-3.603	0.827	0	$d_{35}$	1.121	0.414	0.007	$\ell_{35}$	0.007	0.001	0
$f_{36}$	-0.1	0.596	0.87	$d_{36}$	-0.376	0.231	0.1	$\ell_{36}$	0.001	0.001	0.59

### 6.3. Estimación de parámetros del modelo RDINA con covariable afectando solo atributos

La tabla 6.11 presenta las estimaciones de los parámetros a nivel de atributos de los vectores **b** y **h**. En el caso del vector **b**, las estimaciones de los parámetros  $b_2$ ,  $b_3$ ,  $b_4$  y  $b_7$  fueron significativas con un valor  $p < 0.05$ , lo cual quiere decir, que los examinados con bajos

niveles de Lectura tienden a tener un mayor dominio de las habilidades **R<sub>2</sub>**, **R<sub>3</sub>**, **R<sub>4</sub>** y **R<sub>7</sub>**. En el vector **h**, las estimaciones de los parámetros  $h_2$ ,  $h_3$ ,  $h_4$ ,  $h_5$ ,  $h_6$  y  $h_7$  tuvieron un valor  $p < 0.05$  con valores positivos, por tanto, el puntaje en Lectura tiene un efecto positivo en la probabilidad de que un examinado posea la habilidades **R<sub>2</sub>**, **R<sub>3</sub>**, **R<sub>4</sub>**, **R<sub>5</sub>**, **R<sub>6</sub>** y **R<sub>7</sub>**, que están referidas a la identificación del significado de la palabra, conocimiento de las reglas morfo-sintácticas, las reglas ortográficas, las reglas de puntuación, reconocimiento de la organización textual y comprensión del texto de manera global, respectivamente.

**Tabla 6.11:** Estimación de los vectores de parámetros **b** y **h** del modelo RDINA con covariable afectando solo atributos en la aplicación

Vector de parámetros b				Vector de parámetros h			
Param.	Est.	s.e.	p-valor	Param.	Est.	s.e.	p-valor
$b_1$	-2.342	2.158	0.28	$h_1$	0.006	0.004	0.065
$b_2$	-7.121	2.085	0.001	$h_2$	0.012	0.003	0
$b_3$	-8.586	2.036	0	$h_3$	0.017	0.004	0
$b_4$	-11.981	1.897	0	$h_4$	0.022	0.003	0
$b_5$	-9.327	5.588	0.095	$h_5$	0.021	0.01	0.039
$b_6$	-10.808	6.481	0.095	$h_6$	0.024	0.012	0.04
$b_7$	-5.481	2.309	0.018	$h_7$	0.011	0.004	0.004
$b_8$	-6.134	6.987	0.38	$h_8$	0.015	0.012	0.21

La tabla 6.12 presenta las estimaciones de los parámetros a nivel de ítems de los vectores **f** y **d**. Los parámetros que tuvieron un valor  $p < 0.05$  fueron 23 en el caso del vector **f**, con valores positivos y negativos, y en lo que respecta al vector **d**, fueron 32 parámetros con valores positivos.

#### 6.4. Estimación de parámetros del modelo RDINA sin covariable

La tabla 6.13 presenta las estimaciones de los parámetros a nivel de atributos del vector **b**. Como puede apreciarse, las estimaciones de los parámetros  $b_1$ ,  $b_4$ ,  $b_5$ ,  $b_6$ ,  $b_7$  y  $b_8$  fueron positivas y significativas con un valor  $p < 0.05$ , lo cual quiere decir, que los examinados tienden a tener un mayor dominio de las habilidades **R<sub>1</sub>**, **R<sub>4</sub>**, **R<sub>5</sub>**, **R<sub>6</sub>**, **R<sub>7</sub>** y **R<sub>8</sub>**, referidas a la interpretación del término/expresión a partir del co-texto, conocimiento de las reglas ortográficas, las reglas de puntuación, reconocimiento de la organización textual, comprensión del texto de manera global y uso/identificación de habilidades discursivas académicas, respectivamente.

La tabla 6.14 presenta las estimaciones de los parámetros a nivel de ítems de los vectores **f** y **d**. Los parámetros que tuvieron un valor  $p < 0.05$  fueron 22 en el caso del vector **f**, con valores positivos y negativos, y en lo que respecta al vector **d**, fueron 31 parámetros con valores positivos.

#### 6.5. Comentarios generales

La tabla 6.15 presenta ítems cuyas estimaciones de sus parámetros  $f$ ,  $d$  y  $\ell$  resultaron significativos en los cuatro modelos en estudio. Como se muestra, las estimaciones de los parámetros de adivinación ( $g$ ) y desliz ( $s$ ) en el caso del modelo RDINA con covariable afectando solo atributos fueron ligeramente más bajos en 3 ítems respecto del RDINA sin covariable.

**Tabla 6.12:** Estimación de los vectores de parámetros  $\mathbf{f}$  y  $\mathbf{d}$  del modelo RDINA con covariable afectando solo atributos en la aplicación

Vector de parámetros $\mathbf{f}$				Vector de parámetros $\mathbf{d}$			
Param.	Est.	s.e.	p-valor	Param.	Est.	s.e.	p-valor
$f_1$	0.943	0.142	0	$d_1$	1.231	0.386	0.001
$f_2$	2.498	0.346	0	$d_2$	4.193	2.927	0.15
$f_3$	-0.062	0.115	0.59	$d_3$	1.314	0.327	0
$f_4$	-0.38	0.168	0.023	$d_4$	1.243	0.221	0
$f_5$	0.266	0.147	0.071	$d_5$	0.647	0.21	0.002
$f_6$	0.057	0.146	0.7	$d_6$	0.562	0.204	0.006
$f_7$	1.272	0.172	0	$d_7$	1.672	0.384	0
$f_8$	0.652	0.182	0	$d_8$	1.875	0.294	0
$f_9$	1.108	0.171	0	$d_9$	1.417	0.297	0
$f_{10}$	0.41	0.135	0.002	$d_{10}$	2.104	0.306	0
$f_{11}$	-0.162	0.156	0.3	$d_{11}$	1.156	0.215	0
$f_{12}$	1.144	0.168	0	$d_{12}$	1.486	0.32	0
$f_{13}$	-0.574	0.151	0	$d_{13}$	0.105	0.207	0.61
$f_{14}$	1.339	0.178	0	$d_{14}$	1.986	0.44	0
$f_{15}$	0.067	0.116	0.56	$d_{15}$	1.298	0.324	0
$f_{16}$	1.632	0.191	0	$d_{16}$	1.557	0.383	0
$f_{17}$	-0.293	0.143	0.041	$d_{17}$	2.598	0.337	0
$f_{18}$	0.554	0.137	0	$d_{18}$	2.184	0.347	0
$f_{19}$	0.129	0.27	0.63	$d_{19}$	2.261	0.371	0
$f_{20}$	1.218	0.233	0	$d_{20}$	2.516	0.582	0
$f_{21}$	-0.242	0.186	0.19	$d_{21}$	2.606	0.65	0
$f_{22}$	-0.231	0.199	0.25	$d_{22}$	1.225	0.287	0
$f_{23}$	-1.276	0.227	0	$d_{23}$	0.425	0.292	0.15
$f_{24}$	-1.852	0.235	0	$d_{24}$	1.666	0.273	0
$f_{25}$	0.119	0.178	0.5	$d_{25}$	1.796	0.264	0
$f_{26}$	0.887	0.205	0	$d_{26}$	0.622	0.28	0.026
$f_{27}$	0.978	0.171	0	$d_{27}$	2.422	0.442	0
$f_{28}$	0.009	0.164	0.96	$d_{28}$	1.692	0.248	0
$f_{29}$	-0.66	0.152	0	$d_{29}$	2.194	0.244	0
$f_{30}$	-1.559	0.262	0	$d_{30}$	1.052	0.317	0.001
$f_{31}$	0.877	0.218	0	$d_{31}$	0.78	0.292	0.008
$f_{32}$	1.126	0.184	0	$d_{32}$	1.161	0.302	0
$f_{33}$	0.308	0.179	0.085	$d_{33}$	0.997	0.251	0
$f_{34}$	1.113	0.146	0	$d_{34}$	2.076	0.379	0
$f_{35}$	-0.447	0.362	0.22	$d_{35}$	2.233	0.394	0
$f_{36}$	-0.132	0.144	0.36	$d_{36}$	0.342	0.2	0.087

**Tabla 6.13:** Estimación del vector de parámetros  $\mathbf{b}$  del modelo RDINA sin covariable en la aplicación

Param.	Est.	s.e.	p-valor
$b_1$	1.281	0.304	0
$b_2$	0.439	0.465	0.34
$b_3$	0.079	0.487	0.87
$b_4$	0.709	0.238	0.003
$b_5$	3.506	1.265	0.006
$b_6$	3.405	1.282	0.008
$b_7$	0.778	0.248	0.002
$b_8$	2.665	1.047	0.011

Cuando comparamos los otros dos modelos, el RDINA con covariable afectando ítems y atributos respecto del RDINA con covariable afectando solo ítems, los parámetros de adivinación ( $g$ ) fueron ligeramente más bajos mientras que los parámetros de desliz ( $s$ ) disminuyó en dos casos y en los otros dos aumentó ligeramente. Por tanto, debido al efecto de la covariable Lectura, disminuyó marginalmente la probabilidad de que los examinados que no posean todos los atributos requeridos para dichos ítems de Redacción, puedan adivinar y responder correctamente al ítem. De igual manera, debido al efecto de la misma covariable, disminuyó en dos casos la probabilidad de que los examinados que posean todos los atributos requeridos

**Tabla 6.14:** Estimación de los vectores de parámetros **f** y **d** del modelo RDINA sin covariable en la aplicación

Vector de parámetros f				Vector de parámetros d			
Param.	Est.	s.e.	p-valor	Param.	Est.	s.e.	p-valor
$f_1$	0.87	0.185	0	$d_1$	1.15	0.393	0.003
$f_2$	2.948	0.349	0	$d_2$	3.493	5.11	0.49
$f_3$	-0.16	0.137	0.24	$d_3$	2.331	1.018	0.022
$f_4$	-0.281	0.189	0.14	$d_4$	1.059	0.247	0
$f_5$	0.263	0.152	0.083	$d_5$	0.651	0.216	0.003
$f_6$	0.06	0.153	0.69	$d_6$	0.556	0.214	0.01
$f_7$	1.227	0.179	0	$d_7$	1.82	0.397	0
$f_8$	0.731	0.207	0	$d_8$	1.602	0.318	0
$f_9$	1.153	0.175	0	$d_9$	1.296	0.302	0
$f_{10}$	0.459	0.139	0.001	$d_{10}$	1.954	0.304	0
$f_{11}$	-0.104	0.161	0.52	$d_{11}$	1.053	0.22	0
$f_{12}$	1.107	0.174	0	$d_{12}$	1.59	0.342	0
$f_{13}$	-0.549	0.152	0	$d_{13}$	0.064	0.21	0.76
$f_{14}$	1.292	0.185	0	$d_{14}$	2.182	0.449	
$f_{15}$	0.052	0.132	0.7	$d_{15}$	1.212	0.346	0
$f_{16}$	1.652	0.197	0	$d_{16}$	1.492	0.388	0
$f_{17}$	-0.404	0.161	0.012	$d_{17}$	3.115	0.486	0
$f_{18}$	0.641	0.14	0	$d_{18}$	1.863	0.316	0
$f_{19}$	0.163	0.251	0.52	$d_{19}$	2.263	0.375	0
$f_{20}$	1.311	0.227	0	$d_{20}$	2.252	0.514	0
$f_{21}$	-0.366	0.258	0.16	$d_{21}$	2.216	0.625	0
$f_{22}$	-0.184	0.235	0.43	$d_{22}$	1.64	0.441	0
$f_{23}$	-1.471	0.229	0	$d_{23}$	0.869	0.308	0.005
$f_{24}$	-1.958	0.261	0	$d_{24}$	1.8	0.292	0
$f_{25}$	0.005	0.223	0.98	$d_{25}$	1.97	0.32	0
$f_{26}$	1.126	0.194	0	$d_{26}$	0.33	0.328	0.31
$f_{27}$	0.957	0.171	0	$d_{27}$	2.522	0.559	0
$f_{28}$	0.098	0.158	0.54	$d_{28}$	1.501	0.249	0
$f_{29}$	-0.722	0.17	0	$d_{29}$	2.374	0.266	0
$f_{30}$	-1.584	0.319	0	$d_{30}$	1.356	0.349	0
$f_{31}$	1.225	0.193	0	$d_{31}$	0.296	0.33	0.37
$f_{32}$	1.261	0.182	0	$d_{32}$	0.86	0.288	0.003
$f_{33}$	0.222	0.162	0.17	$d_{33}$	1.159	0.239	0
$f_{34}$	1.107	0.151	0	$d_{34}$	2.163	0.4	0
$f_{35}$	0.205	0.275	0.46	$d_{35}$	1.581	0.413	0
$f_{36}$	-0.123	0.147	0.4	$d_{36}$	0.327	0.207	0.11

para dichos ítems de Redacción respondan incorrectamente al ítem , es decir, que tengan un desliz y en los otros dos casos dicha probabilidad aumentó marginalmente. Asimismo, puede notarse, que los parámetros de adivinación ( $g$ ) en la tabla 6.15 de los modelos RDINA con covariable afectando ítems y atributos y el RDINA con covariable afectando solo ítems, fueron mayores por efecto de la covariable Lectura, a los de adivinación ( $g$ ) de los modelos RDINA con covariable afectando solo atributos y RDINA sin covariable.

**Tabla 6.15:** Estimación de los parámetros de adivinación ( $g$ ) y desliz ( $s$ ) para un  $Z = 574,6$ , que equivale a 28 ítems respondidos correctamente en promedio por los examinados en la prueba de Lectura.

Item	Covariable afectando ítems y atributos		Covariable afectando ítems		Covariable afectando atributos		Sin covariable	
N°	Est. $g$	Est. $s$	Est. $g$	Est. $s$	Est. $g$	Est. $s$	Est. $g$	Est. $s$
8	0.769	0.098	0.791	0.106	0.658	0.074	0.675	0.088
10	0.661	0.087	0.680	0.086	0.601	0.075	0.613	0.082
17	0.437	0.047	0.483	0.054	0.427	0.091	0.400	0.062
18	0.712	0.089	0.724	0.085	0.635	0.061	0.655	0.076

Usando las ecuaciones (3.6) y (3.11) calculamos  $p(\alpha_{ck})$ , la probabilidad de que un examinado en la clase  $c$  posea la habilidad  $k$  para aquellos parámetros que resultaron significativos.



En el caso del modelo RDINA con covariable afectando ítems y atributos las probabilidades para las habilidades  $\mathbf{R}_3$  y  $\mathbf{R}_4$  en promedio fueron 0.4168 y 0.6108, respectivamente. En lo que respecta al RDINA con covariable afectando solo ítems las probabilidades para las habilidades  $\mathbf{R}_1$  y  $\mathbf{R}_6$  fueron 0.8248 y 0.9286, respectivamente. En el RDINA con covariable afectando solo atributos, las probabilidades para las habilidades  $\mathbf{R}_2$ ,  $\mathbf{R}_3$ ,  $\mathbf{R}_4$  y  $\mathbf{R}_7$  en promedio fueron 0.4792, 0.6874, 0.6515 y 0.7330, respectivamente. En el caso del RDINA sin covariable, las probabilidades para las habilidades  $\mathbf{R}_1$ ,  $\mathbf{R}_4$ ,  $\mathbf{R}_5$ ,  $\mathbf{R}_6$ ,  $\mathbf{R}_7$  y  $\mathbf{R}_8$  fueron 0.7826, 0.6701, 0.9709, 0.9679, 0.6852 y 0.9349, respectivamente.

Puede notarse los cambios por efecto de la covariable en la probabilidad de que un examinado en una clase  $c$  posea el atributo  $\mathbf{R}_4$ , en la cual se obtuvo 0.6108 cuando la covariable afectaba los ítems y atributos y 0.6515 cuando afectaba solo a los atributos. Por tanto, el efecto de la covariable Lectura, además de impactar en el dominio del atributo  $\mathbf{R}_4$ , afectan a su vez, a la probabilidad de pertenencia del mismo examinado a las clases que tengan dentro de su estado de conocimiento el atributo  $\mathbf{R}_4$ .

Lo mismo puede concluirse con respecto a los cambios, en la probabilidad de que un examinado en una clase  $c$  posea el atributo  $\mathbf{R}_7$  en la cual se obtuvo 0.7330 cuando la covariable afectaba solo los atributos, y 0.6852 cuando la covariable no le afectaba. Este cambio, también afecta a su vez, a la probabilidad de pertenencia del mismo examinado a las clases que tengan dentro de su estado de conocimiento el atributo  $\mathbf{R}_7$ .

Por otro lado, en las tablas 6.16 y 6.17 muestra para algunos ítems que se encuentran en el Apéndice D, los cambios en la estimación de los parámetros de adivinación ( $g$ ) y desliz ( $s$ ) del modelo RDINA con covariable afectando ítems y atributos y del RDINA con covariable afectando solo ítems, cuando se incrementa el puntaje en Lectura de  $Z = 574,6$  que corresponde a 28 ítems respondidos, a  $Z = 594,7$  que corresponde a 29 ítems respondidos. Dicha equivalencia corresponde a la tabla que maneja la universidad privada para su prueba de admisión, que es objeto de estudio en la aplicación. Asimismo, se ha seleccionado como ejemplo  $Z = 574,6$  (28 ítems) que corresponde al promedio de respuestas de los examinados.

En la tabla 6.16, del modelo RDINA con covariable afectando ítems y atributos, por ejemplo, si analizamos el ítem 18 de Redacción, donde se pregunta si hay algún error ortográfico en un enunciado, se aprecia que el parámetro de adivinación ( $g$ ) se incrementó en 4.7%, cuando el puntaje en Lectura aumentó de  $Z = 574,6$  a  $Z = 594,7$ , mientras que el parámetro de desliz ( $s$ ) disminuyó en 14.5%. Asimismo, en la tabla 6.17 del modelo RDINA con covariable afectando solo ítems, si analizamos el ítem 28 de Redacción, donde se pide ordenar unas oraciones de forma que construya un texto coherente que responda al título propuesto, se aprecia que el parámetro de adivinación ( $g$ ) se incrementó en 3.5%, cuando el puntaje en Lectura aumentó de  $Z = 574,6$  a  $Z = 594,7$ , mientras que el parámetro de desliz ( $s$ ) disminuyó en 10.6%.

Finalmente, en las tablas 6.18 y 6.19 muestra para algunas habilidades, los cambios en la probabilidad que un examinado posea la habilidad  $k$  ( $p(\alpha_{ck}|Z_i = z)$ ) del modelo RDINA con covariable afectando ítems y atributos y del RDINA con covariable afectando solo atributos, cuando se incrementa el puntaje en Lectura de  $Z = 574,6$  a  $Z = 594,7$ .

En la tabla 6.18, del modelo RDINA con covariable afectando ítems y atributos, por ejemplo, si analizamos la habilidad 3 referido al conocimiento de reglas morfo-sintácticas, se aprecia que la probabilidad que un examinado posea dicha habilidad se incrementó en 9.1%, cuando el puntaje en Lectura aumentó de  $Z = 574,6$  a  $Z = 594,7$ . Asimismo, en la tabla 6.19, del modelo RDINA con covariable afectando solo atributos, si analizamos la habilidad

**Tabla 6.16:** Estimación de los parámetros de adivinación ( $g$ ) y desliz ( $s$ ) del modelo RDINA con covariable afectando ítems y atributos, cuando se incrementa el puntaje en Lectura de  $Z = 574,6$  a  $Z = 594,7$ .

Item	Estimaciones para $Z = 574,6$ (28 ítems)		Estimaciones para $Z = 594,7$ (29 ítems)		% Variación	
Nº	$g$	$s$	$g$	$s$	% $g$	% $s$
4	0.484	0.335	0.507	0.315	4.7 %	-5.9 %
11	0.521	0.305	0.551	0.280	5.9 %	-8.3 %
18	0.712	0.089	0.746	0.076	4.7 %	-14.5 %
28	0.574	0.188	0.614	0.164	6.9 %	-12.7 %
35	0.622	0.172	0.647	0.157	4.1 %	-8.7 %

**Tabla 6.17:** Estimación de los parámetros de adivinación ( $g$ ) y desliz ( $s$ ) del modelo RDINA con covariable afectando solo ítems, cuando se incrementa el puntaje en Lectura de  $Z = 574,6$  a  $Z = 594,7$ .

Item	Estimaciones para $Z = 574,6$ (28 ítems)		Estimaciones para $Z = 594,7$ (29 ítems)		% Variación	
Nº	$g$	$s$	$g$	$s$	% $g$	% $s$
4	0.503	0.332	0.532	0.307	5.9 %	-7.8 %
11	0.695	0.496	0.724	0.462	4.0 %	-6.9 %
18	0.724	0.085	0.765	0.070	5.6 %	-18.0 %
28	0.801	0.421	0.829	0.377	3.5 %	-10.6 %
35	0.631	0.160	0.664	0.142	5.2 %	-11.6 %

7 referido a la comprensión del texto de manera global, se aprecia que la probabilidad que un examinado posea dicha habilidad se incrementó en 6.4 %, cuando el puntaje en Lectura aumentó de  $Z = 574,6$  a  $Z = 594,7$ .

**Tabla 6.18:** Estimación de la probabilidad de poseer la habilidad  $k$  ( $p(\alpha_{ck}|Z_i = z)$ ) del modelo RDINA con covariable afectando ítems y atributos, cuando se incrementa el puntaje en Lectura de  $Z = 574,6$  a  $Z = 594,7$ .

Habilidad	Estimación $p(\alpha_{ck} Z_i = z)$ para $Z = 574,6$ (28 ítems)	Estimación $p(\alpha_{ck} Z_i = z)$ para $Z = 594,7$ (29 ítems)	% Variación
3	0.440	0.480	9.1 %
4	0.566	0.615	8.6 %

**Tabla 6.19:** Estimación de la probabilidad de poseer la habilidad  $k$  ( $p(\alpha_{ck}|Z_i = z)$ ) del modelo RDINA con covariable afectando solo atributos, cuando se incrementa el puntaje en Lectura de  $Z = 574,6$  a  $Z = 594,7$ .

Habilidad	Estimación $p(\alpha_{ck} Z_i = z)$ para $Z = 574,6$ (28 ítems)	Estimación $p(\alpha_{ck} Z_i = z)$ para $Z = 594,7$ (29 ítems)	% Variación
2	0.444	0.504	13.5 %
3	0.765	0.821	7.3 %
4	0.659	0.751	13.9 %
7	0.698	0.743	6.4 %

## Capítulo 7

# Conclusiones y sugerencias

### 7.1. Conclusiones

- En la presente tesis se desarrolló el modelo RDINA y su extensión con el uso de una covariable bajo el enfoque clásico, para lo cual se realizó un estudio de simulación con información del TIMMS, buscando evaluar la recuperación de los parámetros y luego aplicarlo al campo educacional, específicamente a la prueba de admisión de una universidad.
- Se realizaron estudios de simulación bajo 4 modelos: RDINA con covariable afectando ítems y atributos, RDINA con covariable afectando solo ítems, RDINA con covariable afectando solo atributos y RDINA sin covariable. A partir de los resultados obtenidos, se identificaron las siguientes conclusiones:
  - El sesgo porcentual de los vectores de parámetros **b**, **h**, **f**, **d** y **ℓ** fueron similares y en su mayoría menores o iguales al 5 %. En el RDINA con covariable afectando ítems y atributos el 82 % de los 89 parámetros obtuvo un sesgo porcentual igual o menor al 5 %, en el RDINA con covariable afectando solo ítems el 95 % de los 82 parámetros obtuvo un sesgo porcentual igual o menor al 5 %, en el RDINA con covariable afectando solo atributos el 92 % de los 64 parámetros obtuvo un sesgo porcentual igual o menor al 5 % y en el RDINA sin covariable el 91 % de los 57 parámetros obtuvo un sesgo porcentual igual o menor al 5 %.
  - Los parámetros del vector **b** de los modelos RDINA con covariable afectando atributos e ítems y RDINA con covariable afectando solo atributos presentaron en algunos casos mayor sesgo respecto a los modelos RDINA con covariable afectando solo ítems y RDINA sin covariable.
  - Los parámetros del vector **h** del modelo RDINA con covariable afectando atributos e ítems presentaron en varios casos menor sesgo respecto al modelo RDINA con covariable afectando solo atributos.
  - Los sesgos de los parámetros del vector **f** en los 4 modelos presentaron resultados similares y casi la totalidad de ellos menores al 5 %.
  - Los parámetros del vector **d** de los modelos RDINA con covariable afectando ítems y atributos y RDINA con covariable afectando solo ítems presentaron en algunos casos mayor sesgo respecto a los modelos RDINA con covariable afectando solo atributos y RDINA sin covariable.

- Los parámetros del vector  $\ell$  del modelo RDINA con covariable afectando atributos e ítems presentaron en general mayor sesgo respecto al modelo RDINA con covariable afectando solo ítems.
- En conclusión, los resultados de simulación mostraron estabilidad al momento de recuperar los parámetros bajo los 4 modelos propuestos dado que su sesgo en la gran mayoría estuvo por debajo del 5%, pero particularmente aquellos modelos con menor cantidad de parámetros, es decir, RDINA con covariable afectando solo ítems (82 parámetros), RDINA con covariable afectando solo atributos (64 parámetros) y RDINA sin covariable (57 parámetros) tuvieron un mejor desempeño pues más del 90 % de sus parámetros obtuvieron un sesgo porcentual igual o menor al 5 %, mientras que en el modelo RDINA con covariable afectando ítems y atributos (89 parámetros) el 82 % de sus parámetros obtuvieron un sesgo porcentual igual o menor al 5 %. Los resultados obtenidos era lo que se esperaba, pues guardan concordancia con el estudio de simulación realizado por Yoon Soo Park (2014).
- Los cuatro estudios de simulación se realizaron simulando las matrices de respuestas dicotómicas usando el software R, luego se usó el software Latent GOLD 5.1 para la estimación, el cual utiliza los algoritmos de Esperanza-Maximización (EM) y de Newton-Raphson. El tiempo de demora máximo para la recuperación de parámetros en el Latent fue de 15 minutos por réplica en una computadora con procesador Intel Core i7, memoria RAM 16.0 GB para el modelo RDINA con covariable afectando atributos e ítems que tiene la mayor cantidad de parámetros (89).
- Respecto a la aplicación en una prueba de admisión de una universidad privada, se observan las siguientes conclusiones:
  - Se ha verificado la identificabilidad de los modelos en la aplicación, para lo cual se usaron 100 conjuntos de valores iniciales en el algoritmo, a diferencia del estudio de simulación en el que se usaron solo 20, permitiendo obtener las mismas estimaciones de los parámetros finales.
  - En las estimaciones de los parámetros del RDINA con covariable afectando atributos e ítems, 21 parámetros de los 36 del vector  $\ell$  tuvieron un efecto significativo con valores positivos. Por lo tanto, la covariable Lectura tuvo influencia positiva en la probabilidad de responder correctamente a estos 21 ítems de Redacción. Asimismo, 3 parámetros de los 8 del vector  $\mathbf{h}$  tuvieron un efecto significativo con valores positivos, es decir, la covariable Lectura también tuvo influencia positiva en la probabilidad que un examinado posea las habilidades asociadas a estos 3 parámetros.
  - En las estimaciones de los parámetros del RDINA con covariable afectando solo ítems, 30 parámetros de los 36 del vector  $\ell$  tuvieron un efecto significativo con valores positivos. Por lo tanto, la covariable Lectura tuvo influencia positiva en la probabilidad de responder correctamente a estos 30 ítems de Redacción.
  - En las estimaciones de los parámetros del RDINA con covariable afectando solo atributos, 6 parámetros de los 8 del vector  $\mathbf{h}$  tuvieron un efecto significativo con valores positivos. Por lo tanto, la covariable Lectura tuvo influencia positiva en la probabilidad que un examinado posea las habilidades asociadas a estos 6 parámetros.
- De lo anterior, podemos concluir con respecto a la aplicación, que en el modelo RDINA con covariable afectando solo ítems, la covariable Lectura tuvo influencia positiva en la probabilidad de responder correctamente un ítem en una mayor cantidad de preguntas



de Redacción en comparación al modelo RDINA con covariable afectando atributos e ítems. De igual manera, en el modelo RDINA con covariable afectando solo atributos, la covariable Lectura tuvo influencia positiva en la probabilidad de poseer una habilidad en una mayor cantidad de ellas en comparación al modelo RDINA con covariable afectando atributos e ítems.

- En la aplicación se obtuvieron los resultados esperados referente al efecto de la covariable Lectura en los ítems de la prueba de Redacción, pues como se apreció en las tablas 6.16 y 6.17, los parámetros de adivinación ( $g$ ) tuvieron un incremento porcentual, es decir, aumentó la probabilidad de un examinado en responder correctamente un ítem de Redacción a pesar de no tener todas las habilidades requeridas para dicho ítem, cuando se aumentó el puntaje en Lectura, mientras que los parámetros de desliz ( $s$ ) tuvieron una disminución porcentual, es decir, se redujo la probabilidad de un examinado en responder incorrectamente un ítem de Redacción, a pesar de poseer todas las habilidades requeridas para dicho ítem, cuando se aumentó el puntaje en Lectura.
- Asimismo, se obtuvieron los resultados esperados referente al efecto de la covariable Lectura en las habilidades, pues como se apreció en las tablas 6.18 y 6.19, la probabilidad que un examinado posea una determinada habilidad tuvo un incremento porcentual, cuando se aumentó el puntaje en Lectura.
- Finalmente, cabe señalar la importancia en la construcción de la matriz  $Q$ , que debe ser realizado por un equipo de especialistas en las competencias que se están evaluando. En el caso de la aplicación se construyó por un especialista en base a una prueba de admisión enfocada más en el dominio de áreas temáticas que en el dominio de habilidades, lo cual puede haber influido en la significancia de los parámetros estimados.

## 7.2. Sugerencias para investigaciones futuras

A continuación, se indican sugerencias para investigaciones futuras que puedan surgir a partir del estudio realizado en el presente trabajo.

- Si bien los resultados del estudio de simulación para un modelo RDINA con una covariable donde hay 7 atributos, mostraron buen desempeño al momento de recuperar los parámetros bajo los cuatro modelos propuestos con un tamaño de muestra de 500, es posible que dicho tamaño requerido deba aumentar a 1000 o 2000 examinados si el número de atributos aumenta a fin de que el sesgo en las estimaciones de los parámetros se mantenga por debajo del 5 %.
- Si bien el estudio de simulación se realizó con una longitud de prueba de 25 ítems, también debería estudiarse para diversas longitudes de prueba, atributos y especificación de atributos indicados en la matriz  $Q$ . Asimismo, debe considerarse en futuras simulaciones estudios que incluyan un mayor número de covariables y de diferentes tipos como dicotómicas u ordinales, también una combinación de covariables continuas y discretas que incluyan características demográficas.
- La extensión de la covariable realizado al modelo DINA se hizo bajo el supuesto de independencia condicional para las respuestas a las clases latentes de los examinados y también en la estructura de atributos. Por otro lado, la especificación del CDM es independiente a la especificación del atributo; por lo que puede usarse junto con cualquier especificación de atributo, como saturación (modelo completo), independencia



o independencia condicional. En el presente estudio se examinó el uso de especificaciones de atributos restringidos (independencia e independencia condicional); sin embargo, una estructura de atributos saturada sin restricciones representa el modelo más general para la distribución de atributos. Por tanto, pueden requerirse estudios futuros que examinen las comparaciones de especificaciones de atributos restringidos y no restringidos, en el contexto de las extensiones de covariables.



## Apéndice A

# Algoritmo Esperanza-Maximización (EM)

### A.1. Cálculo del paso Maximización(M)

Dado  $\mathbf{L}_c(\boldsymbol{\theta}; \mathbf{y}, \mathbf{z}, \mathbf{s})$  el espacio de parámetros de  $\boldsymbol{\theta}$  será conveniente separar  $(J+1)$  vectores:

- Un vector corresponde a las habilidades  $\boldsymbol{\beta}_k = (\mathbf{b}_k, \mathbf{h}_k)$  con  $\mathbf{b}_k = (b_1, \dots, b_K)$  y  $\mathbf{h}_k = (h_1, \dots, h_K)$ ,
- Los otros  $J$  vectores corresponden a los ítems  $\mathbf{w}_j = (\mathbf{f}_j, \mathbf{d}_j, \ell_j)$  con  $\mathbf{f}_j = (f_1, \dots, f_J)$ ,  $\mathbf{d}_j = (d_1, \dots, d_J)$  y  $\ell_j = (\ell_1, \dots, \ell_J)$ , para  $j = 1, \dots, J$ .

Luego con la finalidad de realizar la maximización del algoritmo EM, es decir, encontrar  $\boldsymbol{\theta}$  que maximiza  $Q(\boldsymbol{\theta}|\boldsymbol{\theta}^{(p)})$  implementaremos por separado cada vector, ahorrando de esta manera tiempo computacional. Por tanto el paso **M** puede escribirse como:

Encontrar  $\boldsymbol{\beta}_k$  que maximice:

$$\begin{aligned} Q_{\boldsymbol{\beta}_k}(\boldsymbol{\beta}_k|\boldsymbol{\theta}^{(p)}) &= E \left( \sum_{i=1}^N \sum_{c=1}^C s_{ic} \log P(\boldsymbol{\alpha}_c|Z_i = z) \right) \\ &= \sum_{i=1}^N \sum_{c=1}^C E(s_{ic}) \log P(\boldsymbol{\alpha}_c|Z_i = z) \\ &= \sum_{i=1}^N \sum_{c=1}^C \phi_{ic}(\boldsymbol{\theta}^{(p)}) \sum_{k=1}^K \log \left( \frac{e^{b_k + h_k z_i}}{1 + e^{b_k + h_k z_i}} \right), \end{aligned} \quad (\text{A.1})$$

y encontrar  $\mathbf{w}_j$ ,  $j = 1, \dots, J$ , que maximice

$$\begin{aligned}
\mathbf{Q}_{\mathbf{w}_j}(\mathbf{w}_j|\boldsymbol{\theta}^{(p)}) &= \mathbb{E} \left( \sum_{i=1}^N \sum_{c=1}^C s_{ic} \log P(Y_{ij} = y_{ij}|\boldsymbol{\alpha}_c, Z_i = z) \right) \\
&= \sum_{i=1}^N \sum_{c=1}^C \mathbb{E}(s_{ic}) \log P(Y_{ij} = y_{ij}|\boldsymbol{\alpha}_c, Z_i = z) \\
&= \sum_{i=1}^N \sum_{c=1}^C \phi_{ic}(\boldsymbol{\theta}^{(p)}) \log \left( \left( \frac{e^{f_j + d_j \eta_{cj} + \ell_j z_i}}{1 + e^{f_j + d_j \eta_{cj} + \ell_j z_i}} \right)^{y_{ij}} \left( \frac{1}{1 + e^{f_j + d_j \eta_{cj} + \ell_j z_i}} \right)^{1-y_{ij}} \right), \quad (\text{A.2})
\end{aligned}$$

donde

$$\begin{aligned}
\phi_{ic}(\boldsymbol{\theta}^{(p)}) &= E(s_{ic}|\mathbf{Y}_i = \mathbf{y}_i, \boldsymbol{\theta}^{(p)}, z) = \\
&= \frac{\prod_{k=1}^K p(\alpha_{ck}) \prod_{j=1}^J P(Y_{ij} = y_{ij}|\boldsymbol{\alpha}_c, z)}{\sum_{l=1}^C \prod_{k=1}^K p(\alpha_{lk}) \prod_{j=1}^J P(Y_{ij} = y_{ij}|\boldsymbol{\alpha}_l, z)} = \\
&= \frac{\prod_{k=1}^K \left( \frac{e^{b_k^{(p)} + h_k^{(p)} z_i}}{1 + e^{b_k^{(p)} + h_k^{(p)} z_i}} \right) \prod_{j=1}^J \left( \left( \frac{e^{f_j^{(p)} + d_j^{(p)} \eta_{cj} + \ell_j^{(p)} z_i}}{1 + e^{f_j^{(p)} + d_j^{(p)} \eta_{cj} + \ell_j^{(p)} z_i}} \right)^{y_{ij}} \left( \frac{1}{1 + e^{f_j^{(p)} + d_j^{(p)} \eta_{cj} + \ell_j^{(p)} z_i}} \right)^{1-y_{ij}} \right)}{\sum_{\ell=1}^C \prod_{k=1}^K \left( \frac{e^{b_k^{(p)} + h_k^{(p)} z_i}}{1 + e^{b_k^{(p)} + h_k^{(p)} z_i}} \right) \prod_{j=1}^J \left( \left( \frac{e^{f_j^{(p)} + d_j^{(p)} \eta_{lj} + \ell_j^{(p)} z_i}}{1 + e^{f_j^{(p)} + d_j^{(p)} \eta_{lj} + \ell_j^{(p)} z_i}} \right)^{y_{ij}} \left( \frac{1}{1 + e^{f_j^{(p)} + d_j^{(p)} \eta_{lj} + \ell_j^{(p)} z_i}} \right)^{1-y_{ij}} \right)}, \quad (\text{A.3})
\end{aligned}$$

con  $\ell = 1, \dots, C-1$  y en la cual  $\phi_{ic}(\boldsymbol{\theta}^{(p)})$  es la probabilidad posterior de pertenencia a clase evaluada en  $\boldsymbol{\theta}^{(p)}$ . El proceso de maximización se lleva a cabo utilizando el método de Newton-Raphson de una iteración. Las primeras derivadas en los métodos de Newton para obtener el máximo de las estimaciones de los parámetros de  $\mathbf{Q}_{\beta_k}(\beta_k|\boldsymbol{\theta}^{(p)})$  y  $\mathbf{Q}_{\mathbf{w}_j}(\mathbf{w}_j|\boldsymbol{\theta}^{(p)})$  quedan como sigue:

$$\begin{aligned}
\frac{\partial \mathbf{Q}_{\beta_k}(\beta_k|\boldsymbol{\theta}^{(p)})}{\partial b_k} &= \sum_{i=1}^N (\phi_{ic}^{(p)} - \eta_{ic}) \\
\frac{\partial \mathbf{Q}_{\beta_k}(\beta_k|\boldsymbol{\theta}^{(p)})}{\partial h_k} &= \sum_{i=1}^N (z_i (\phi_{ic}^{(p)} - \eta_{ic})) \\
\frac{\partial \mathbf{Q}_{\mathbf{w}_j}(\mathbf{w}_j|\boldsymbol{\theta}^{(p)})}{\partial f_j} &= \sum_{i=1}^N [\phi_{ic}^{(p)} (y_{ij} - P(Y_{ij} = y_{ij}|\boldsymbol{\alpha}_c, z))] \\
\frac{\partial \mathbf{Q}_{\mathbf{w}_j}(\mathbf{w}_j|\boldsymbol{\theta}^{(p)})}{\partial d_j} &= \sum_{i=1}^N [\eta_{ic} \phi_{ic}^{(p)} (y_{ij} - P(Y_{ij} = y_{ij}|\boldsymbol{\alpha}_c, z))] \\
\frac{\partial \mathbf{Q}_{\mathbf{w}_j}(\mathbf{w}_j|\boldsymbol{\theta}^{(p)})}{\partial \ell_j} &= \sum_{i=1}^N \sum_{c=1}^C [z_i \phi_{ic}^{(p)} (y_{ij} - P(Y_{ij} = y_{ij}|\boldsymbol{\alpha}_c, z))]
\end{aligned}$$

y las segundas derivadas de la siguiente manera:

$$\begin{aligned}
\frac{\partial^2 \mathbf{Q}_{\beta_k}(\beta_k | \boldsymbol{\theta}^{(p)})}{\partial b_k \partial b_{k'}} &= - \sum_{i=1}^N (\delta_{cl} - \eta_{il}) \\
\frac{\partial^2 \mathbf{Q}_{\beta_k}(\beta_k | \boldsymbol{\theta}^{(p)})}{\partial h_k \partial h_{k'}} &= - \sum_{i=1}^N [z_i^2 \eta_{ic} (\delta_{cl} - \eta_{il})] \\
\frac{\partial^2 \mathbf{Q}_{\beta_k}(\beta_k | \boldsymbol{\theta}^{(p)})}{\partial b_k \partial h_k} &= - \sum_{i=1}^N [z_i \eta_{ic} (\delta_{cl} - \eta_{il})] \\
\frac{\partial^2 \mathbf{Q}_{\mathbf{w}_j}(\mathbf{w}_j | \boldsymbol{\theta}^{(p)})}{\partial f_j \partial f_{j'}} &= - \sum_{i=1}^N [\phi_{ic'}^{(p)} \delta_{c'l'} P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_{c'}, z) (\delta_{rs} - P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_{l'}, z))] \\
\frac{\partial^2 \mathbf{Q}_{\mathbf{w}_j}(\mathbf{w}_j | \boldsymbol{\theta}^{(p)})}{\partial d_j \partial d_{j'}} &= - \sum_{i=1}^N [\phi_{ic'}^{(p)} \eta_{ic'}^2 \delta_{c'l'} P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_{c'}, z) (\delta_{rs} - P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_{l'}, z))] \\
\frac{\partial^2 \mathbf{Q}_{\mathbf{w}_j}(\mathbf{w}_j | \boldsymbol{\theta}^{(p)})}{\partial l_j \partial l_{j'}} &= - \sum_{i=1}^N \sum_{c'=1}^C [z_i^2 \phi_{ic'}^{(p)} P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_{c'}, z) (\delta_{rs} - P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_{c'}, z))] \\
\frac{\partial^2 \mathbf{Q}_{\mathbf{w}_j}(\mathbf{w}_j | \boldsymbol{\theta}^{(p)})}{\partial f_j \partial d_j} &= - \sum_{i=1}^N [\phi_{ic'}^{(p)} \eta_{ic'} \delta_{c'l'} P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_{c'}, z) (\delta_{rs} - P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_{l'}, z))] \\
\frac{\partial^2 \mathbf{Q}_{\mathbf{w}_j}(\mathbf{w}_j | \boldsymbol{\theta}^{(p)})}{\partial f_j \partial l_j} &= - \sum_{i=1}^N [z_i \phi_{ic'}^{(p)} P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_{c'}, z) (\delta_{rs} - P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_{c'}, z))] \\
\frac{\partial^2 \mathbf{Q}_{\mathbf{w}_j}(\mathbf{w}_j | \boldsymbol{\theta}^{(p)})}{\partial d_j \partial l_j} &= - \sum_{i=1}^N [z_i \eta_{ic'} \phi_{ic'}^{(p)} P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_{c'}, z) (\delta_{rs} - P(Y_{ij} = y_{ij} | \boldsymbol{\alpha}_{c'}, z))]
\end{aligned}$$

con  $c', l' = 1, \dots, C$ ;  $\delta_{cl} = I(c = l)$ ;  $\eta_{ic} = \eta_c$  y  $r, s = 0, 1$

## Apéndice B

# Estudio de Simulación

### B.1. Valores iniciales, N=500, K=7, J=25

Modelo RDINA con covariable afectando ítems y atributos

$f=c(-1.58,-5.82,-3.32,-1.77,-4.36,-2.97,-3.19,-3.17,-3.58,-2.41,-1.63,-1.01,-2.36,-3.06,-2.36,-5.89,-4.66,-3.45,-5.11,-1.94,-5.02,-1.74,-4.58,-3.36,-3.27)$

$d=c(1.51,2.26,2.06,1.44,2.63,2.87,2.54,1.73,2.66,1.56,0.45,1.04,2.05,1.55,1.56,2.53,1.76,1.86,2.19,1.83,1.08,0.77,2.31,1.88,2.81)$

$l=c(0.15,0.21,0.13,0.13,0.23,0.31,0.17,0.15,0.21,0.12,0.05,0.12,0.18,0.13,0.17,0.28,0.15,0.19,0.23,0.11,0.16,0.08,0.24,0.1,0.23)$

$b=c(-5.65,-5.52,-3.06,-0.55,-3.04,-2.77,-3.58)$

$h=c(0.34,0.33,0.24,0.08,0.21,0.18,0.23)$

Modelo RDINA con covariable afectando solo ítems

$f=c(-1.49,-5.25,-3.07,-1.74,-4.19,-2.3,-2.8,-3.04,-3.13,-2.13,-1.56,-0.85,-2.33,-2.94,-2.33,-5.54,-4.55,-3.37,-4.96,-1.89,-4.98,-1.61,-4.36,-3.06,-3.04)$

$d=c(1.41,2.09,1.99,1.3,2.34,3.78,1.48,1.63,2.41,1.46,0.47,0.89,1.61,1.26,1.38,2.26,1.62,1.61,1.97,1.68,1.08,0.51,2.13,1.71,2.41)$

$l=c(0.15,0.19,0.12,0.13,0.23,0.27,0.15,0.14,0.19,0.1,0.05,0.11,0.18,0.13,0.17,0.27,0.15,0.19,0.23,0.11,0.16,0.07,0.23,0.09,0.22)$

$b=c(0.14,0.59,2,0.94,0.47,0.9,0.67)$

Modelo RDINA con covariable afectando solo atributos

$f=c(0.89,-2.24,-0.77,0.21,-0.45,1.53,-0.57,-0.92,-0.27,-0.69,-0.88,0.84,0.52,-1.02,0.21,-1.5,-1.94,-0.41,-1.37,-0.11,-2.7,-0.66,-0.8,-1.56,0.44)$

$d=c(1.48,3.12,1.62,2.06,2.83,5.27,2.23,1.93,2.83,1.72,0.56,1.34,2.08,2.01,2.25,3.04,1.7,2.37,2.77,2.24,1.94,1.11,2.75,1.86,2.75)$

$b=c(-5.67,-5.75,-2.6,-1.5,-3.7,-6.33,-3.36)$



```
h=c(0.34,0.34,0.21,0.14,0.24,0.4,0.22)
```

#### Modelo RDINA sin covariable

```
f=c(0.91,-2.43,-0.94,0.14,-0.42,1.31,-0.71,-0.88,-0.24,-0.79,-0.88,0.82,0.49,-1.08,  
0.29,-1.37,-1.96,-0.41,-1.32,-0.11,-2.58,-0.61,-0.74,-1.62,0.42)
```

```
d=c(1.51,2.58,1.78,1.93,2.75,4.6,1.87,1.99,2.62,1.65,0.53,1.28,2.14,1.91,2.13,2.98,  
1.74,2.35,2.65,2.07,1.81,0.91,2.78,1.91,2.49)
```

```
b=c(0.10,0.63,1.69,1.51,0.45,0.95,0.75)
```

## B.2. Matriz Q de habilidades

```
Q1<-matrix(c(1,0,0,0,0,0,0,0,1,0,0,0,0,0,1,1,0,0,0,0,0,  
1,1,0,0,0,0,0,1,0,1,0,0,0,0,0,0,0,0,1,1,0,  
0,0,0,1,1,1,0,1,0,0,0,1,0,0,0,0,0,0,1,0,0,  
0,0,0,1,1,0,0,1,0,0,1,0,0,0,1,0,0,0,0,0,1,  
1,0,0,0,0,0,1,1,1,0,0,0,0,1,1,0,0,0,0,0,0,  
1,0,0,0,0,0,0,1,0,1,0,0,0,0,1,0,1,0,0,0,0,  
1,0,0,0,0,0,1,1,0,1,0,0,0,1,1,0,1,0,0,0,0,  
0,0,0,0,1,1,0,1,0,0,0,0,0,0,0,0,0,0,1,0,0,  
1,0,0,0,0,0,1),nrow=25,byrow = T)
```

## B.3. Código R - Generación de Bases de Datos, N=500, K=7, J=25

### B.3.1. Modelo RDINA con covariable afectando ítems y atributos

```
##Generación de valores aleatorios de la covariable##
```

```
z<-rnorm(500,17.18,5)
```

```
##Generación de Bases de datos##
```

```
BaseDatos<-function(n){
```

```
Matriz<-matrix(NA,n,25)
```

```
for (i in 1:n){
```

```
prob_alpha<-exp(b+h*z[i])/(1+exp(b+h*z[i]))
```

```
alpha=1
```

```
for(k in 1:7){ alpha[k]=rbinom(1,1,prob_alpha[k]) }
```

```
eta=1
```

```
for(j in 1:25){ eta[j]=alpha[1]^Q1[j,1]*alpha[2]^Q1[j,2]*alpha[3]^Q1[j,3]*  
alpha[4]^Q1[j,4]*alpha[5]^Q1[j,5]*alpha[6]^Q1[j,6]*alpha[7]^Q1[j,7]
```

```
prob_Y<-exp(f+d*eta+l*z[i])/(1+exp(f+d*eta+l*z[i]))
```

```
Y =1
```

```
for(j in 1:25){ Y[j]=rbinom(1,1,prob_Y[j]) }
```

```
Matriz[i,]<-Y
```

```
} Matriz }
```

```
Simulaciones<- lapply(rep(500,100),BaseDatos)
```

### B.3.2. Modelo RDINA con covariable afectando solo ítems

```
##Generación de valores aleatorios de la covariable##

z<-rnorm(500,17.18,5)

##Generación de Bases de datos##

BaseDatos<-function(n){

Matriz<-matrix(NA,n,25)

for (i in 1:n){

  prob_alpha<-exp(b)/(1+exp(b))
  alpha=1
  for(k in 1:7){ alpha[k]=rbinom(1,1,prob_alpha[k]) }
  eta=1
  for(j in 1:25){ eta[j]=alpha[1]^Q1[j,1]*alpha[2]^Q1[j,2]*alpha[3]^Q1[j,3]*
alpha[4]^Q1[j,4]*alpha[5]^Q1[j,5]*alpha[6]^Q1[j,6]*alpha[7]^Q1[j,7]

  prob_Y<-exp(f+d*eta+1*z[i])/(1+exp(f+d*eta+1*z[i]))
  Y =1
  for(j in 1:25){ Y[j]=rbinom(1,1,prob_Y[j]) }

Matriz[i,]<-Y
} Matriz }

Simulaciones<- lapply(rep(500,100),BaseDatos)
```

### B.3.3. Modelo RDINA con covariable afectando solo atributos

```
##Generación de valores aleatorios de la covariable##

z<-rnorm(500,17.18,5)

##Generación de Bases de datos##

BaseDatos<-function(n){

Matriz<-matrix(NA,n,25)

for (i in 1:n){

  prob_alpha<-exp(b+h*z[i])/(1+exp(b+h*z[i]))
  alpha=1
  for(k in 1:7){ alpha[k]=rbinom(1,1,prob_alpha[k]) }
  eta=1
  for(j in 1:25){ eta[j]=alpha[1]^Q1[j,1]*alpha[2]^Q1[j,2]*alpha[3]^Q1[j,3]*
alpha[4]^Q1[j,4]*alpha[5]^Q1[j,5]*alpha[6]^Q1[j,6]*alpha[7]^Q1[j,7]

  prob_Y<-exp(f+d*eta)/(1+exp(f+d*eta))
  Y =1
  for(j in 1:25){ Y[j]=rbinom(1,1,prob_Y[j]) }


```

```
Matriz[i,]<-Y
} Matriz }

Simulaciones<- lapply(rep(500,100),BaseDatos)
```

### B.3.4. Modelo RDINA sin covariable

```
##Generación de Bases de datos##

BaseDatos<-function(n){

Matriz<-matrix(NA,n,25)

for (i in 1:n){

  prob_alpha<-exp(b)/(1+exp(b))
  alpha=1
  for(k in 1:7){ alpha[k]=rbinom(1,1,prob_alpha[k]) }
  eta=1
  for(j in 1:25){ eta[j]=alpha[1]^Q1[j,1]*alpha[2]^Q1[j,2]*alpha[3]^Q1[j,3]*
    alpha[4]^Q1[j,4]*alpha[5]^Q1[j,5]*alpha[6]^Q1[j,6]*alpha[7]^Q1[j,7]

  prob_Y<-exp(f+d*eta)/(1+exp(f+d*eta))
  Y =1
  for(j in 1:25){ Y[j]=rbinom(1,1,prob_Y[j]) }

Matriz[i,]<-Y
} Matriz }

Simulaciones<- lapply(rep(500,100),BaseDatos)
```

## B.4. Código Latent Gold, con K=7, J=25

```
options
  algorithm
    tolerance=1e-008 emtolerance=0.01 emiterations=250 niterations=50 ;
  startvalues
    seed=0 sets=20 tolerance=1e-005 iterations=50;
  bayes
    categorical=1 variances=0 latent=1 poisson=0;
  montecarlo
    seed=0 replicates=500 tolerance=1e-008;
  missing excludeall;
  output
    parameters=first standarderrors probmeans=posterior profile bivariateresiduals
    identification classification iterationdetails;
variables
  dependent Item1 cumlogit, Item2 cumlogit, Item3 cumlogit, Item4 cumlogit,
    Item5 cumlogit, Item6 cumlogit, Item7 cumlogit, Item8 cumlogit,
    Item9 cumlogit, Item10 cumlogit, Item11 cumlogit, Item12 cumlogit,
    Item13 cumlogit, Item14 cumlogit, Item15 cumlogit, Item16 cumlogit,
    Item17 cumlogit, Item18 cumlogit, Item19 cumlogit, Item20 cumlogit;
    Item21 cumlogit, Item22 cumlogit, Item23 cumlogit, Item24 cumlogit,
```

```

Item25 cumlogit;
independent Puntaje_Ciencias;
latent
    a1 ordinal 2 score (0 1), a2 ordinal 2 score (0 1), a3 ordinal 2 score (0 1),
    a4 ordinal 2 score (0 1), a5 ordinal 2 score (0 1), a6 ordinal 2 score (0 1),
    a7 ordinal 2 score (0 1);

```

#### B.4.1. Modelo RDINA con covariable afectando ítems y atributos

equations

```

a1-a7 <- 1 + Puntaje_Ciencias;
Item1 <- 1 + a1 + Puntaje_Ciencias;
Item2 <- 1 + a2 + Puntaje_Ciencias;
Item3 <- 1 + a1 a2 + Puntaje_Ciencias;
Item4 <- 1 + a1 a2 + Puntaje_Ciencias;
Item5 <- 1 + a1 a3 + Puntaje_Ciencias;
Item6 <- 1 + a5 a6 + Puntaje_Ciencias;
Item7 <- 1 + a4 a5 a6 + Puntaje_Ciencias;
Item8 <- 1 + a1 a5 + Puntaje_Ciencias;
Item9 <- 1 + a5 + Puntaje_Ciencias;
Item10 <- 1 + a4 a5 + Puntaje_Ciencias;
Item11 <- 1 + a1 a4 + Puntaje_Ciencias;
Item12 <- 1 + a1 a7 + Puntaje_Ciencias;
Item13 <- 1 + a1 a7 + Puntaje_Ciencias;
Item14 <- 1 + a1 a2 a7 + Puntaje_Ciencias;
Item15 <- 1 + a1 + Puntaje_Ciencias;
Item16 <- 1 + a1 + Puntaje_Ciencias;
Item17 <- 1 + a1 a3 + Puntaje_Ciencias;
Item18 <- 1 + a1 a3 + Puntaje_Ciencias;
Item19 <- 1 + a1 a7 + Puntaje_Ciencias;
Item20 <- 1 + a1 a3 a7 + Puntaje_Ciencias;
Item21 <- 1 + a1 a3 + Puntaje_Ciencias;
Item22 <- 1 + a5 a6 + Puntaje_Ciencias;
Item23 <- 1 + a1 + Puntaje_Ciencias;
Item24 <- 1 + a5 + Puntaje_Ciencias;
Item25 <- 1 + a1 a7 + Puntaje_Ciencias;

```

#### B.4.2. Modelo RDINA con covariable afectando solo ítems

equations

```

a1-a7 <- 1;
Item1 <- 1 + a1 + Puntaje_Ciencias;
Item2 <- 1 + a2 + Puntaje_Ciencias;
Item3 <- 1 + a1 a2 + Puntaje_Ciencias;
Item4 <- 1 + a1 a2 + Puntaje_Ciencias;
Item5 <- 1 + a1 a3 + Puntaje_Ciencias;
Item6 <- 1 + a5 a6 + Puntaje_Ciencias;
Item7 <- 1 + a4 a5 a6 + Puntaje_Ciencias;
Item8 <- 1 + a1 a5 + Puntaje_Ciencias;
Item9 <- 1 + a5 + Puntaje_Ciencias;
Item10 <- 1 + a4 a5 + Puntaje_Ciencias;
Item11 <- 1 + a1 a4 + Puntaje_Ciencias;
Item12 <- 1 + a1 a7 + Puntaje_Ciencias;
Item13 <- 1 + a1 a7 + Puntaje_Ciencias;
Item14 <- 1 + a1 a2 a7 + Puntaje_Ciencias;
Item15 <- 1 + a1 + Puntaje_Ciencias;
Item16 <- 1 + a1 + Puntaje_Ciencias;

```

```

Item17 <- 1 + a1 a3 + Puntaje_Ciencias;
Item18 <- 1 + a1 a3 + Puntaje_Ciencias;
Item19 <- 1 + a1 a7 + Puntaje_Ciencias;
Item20 <- 1 + a1 a3 a7 + Puntaje_Ciencias;
Item21 <- 1 + a1 a3 + Puntaje_Ciencias;
Item22 <- 1 + a5 a6 + Puntaje_Ciencias;
Item23 <- 1 + a1 + Puntaje_Ciencias;
Item24 <- 1 + a5 + Puntaje_Ciencias;
Item25 <- 1 + a1 a7 + Puntaje_Ciencias;

```

#### B.4.3. Modelo RDINA con covariable afectando solo atributos

```

equations
  a1-a7 <- 1 + Puntaje_Ciencias;
  Item1 <- 1 + a1;
  Item2 <- 1 + a2;
  Item3 <- 1 + a1 a2;
  Item4 <- 1 + a1 a2;
  Item5 <- 1 + a1 a3;
  Item6 <- 1 + a5 a6;
  Item7 <- 1 + a4 a5 a6;
  Item8 <- 1 + a1 a5;
  Item9 <- 1 + a5;
  Item10 <- 1 + a4 a5;
  Item11 <- 1 + a1 a4;
  Item12 <- 1 + a1 a7;
  Item13 <- 1 + a1 a7;
  Item14 <- 1 + a1 a2 a7;
  Item15 <- 1 + a1;
  Item16 <- 1 + a1;
  Item17 <- 1 + a1 a3;
  Item18 <- 1 + a1 a3;
  Item19 <- 1 + a1 a7;
  Item20 <- 1 + a1 a3 a7;
  Item21 <- 1 + a1 a3;
  Item22 <- 1 + a5 a6;
  Item23 <- 1 + a1;
  Item24 <- 1 + a5;
  Item25 <- 1 + a1 a7;

```

#### B.4.4. Modelo RDINA sin covariable

```

equations
  a1-a7 <- 1;
  Item1 <- 1 + a1;
  Item2 <- 1 + a2;
  Item3 <- 1 + a1 a2;
  Item4 <- 1 + a1 a2;
  Item5 <- 1 + a1 a3;
  Item6 <- 1 + a5 a6;
  Item7 <- 1 + a4 a5 a6;
  Item8 <- 1 + a1 a5;
  Item9 <- 1 + a5;
  Item10 <- 1 + a4 a5;
  Item11 <- 1 + a1 a4;
  Item12 <- 1 + a1 a7;
  Item13 <- 1 + a1 a7;

```

```
Item14 <- 1 + a1 a2 a7;  
Item15 <- 1 + a1;  
Item16 <- 1 + a1;  
Item17 <- 1 + a1 a3;  
Item18 <- 1 + a1 a3;  
Item19 <- 1 + a1 a7;  
Item20 <- 1 + a1 a3 a7;  
Item21 <- 1 + a1 a3;  
Item22 <- 1 + a5 a6;  
Item23 <- 1 + a1;  
Item24 <- 1 + a5;  
Item25 <- 1 + a1 a7;
```





## Apéndice C

# Aplicación - Códigos Latent Gold

### C.1. Código Latent Gold, con N=727, K=8, J=36

```
options
  algorithm
    tolerance=1e-008 emtolerance=0.01 emiterations=250 nriterations=50 ;
  startvalues
    seed=0 sets=100 tolerance=1e-005 iterations=50;
  bayes
    categorical=1 variances=0 latent=1 poisson=0;
  montecarlo
    seed=0 replicates=500 tolerance=1e-008;
  missing excludeall;
  output
    parameters=first standarderrors probmeans=posterior profile bivariateresiduals
    identification classification iterationdetails;
variables
  dependent Item1 cumlogit, Item2 cumlogit, Item3 cumlogit, Item4 cumlogit,
    Item5 cumlogit, Item6 cumlogit, Item7 cumlogit, Item8 cumlogit, Item9 cumlogit,
    Item10 cumlogit, Item11 cumlogit, Item12 cumlogit, Item13 cumlogit, Item14 cumlogit,
    Item15 cumlogit, Item16 cumlogit, Item17 cumlogit, Item18 cumlogit, Item19 cumlogit,
    Item20 cumlogit, Item21 cumlogit, Item22 cumlogit, Item23 cumlogit, Item24 cumlogit,
    Item25 cumlogit, Item26 cumlogit, Item27 cumlogit, Item28 cumlogit, Item29 cumlogit,
    Item30 cumlogit, Item31 cumlogit, Item32 cumlogit, Item33 cumlogit, Item34 cumlogit,
    Item35 cumlogit, Item36 cumlogit;
  independent Puntaje_Lectura;
  latent
    a1 ordinal 2 score (0 1), a2 ordinal 2 score (0 1), a3 ordinal 2 score (0 1),
    a4 ordinal 2 score (0 1), a5 ordinal 2 score (0 1), a6 ordinal 2 score (0 1),
    a7 ordinal 2 score (0 1), a8 ordinal 2 score (0 1);
```

#### C.1.1. Modelo RDINA con covariable afectando ítems y atributos

```
equations
  a1-a8 <- 1 + Puntaje_Lectura;
  Item1 <- 1 + a1 a5 a8 + Puntaje_Lectura;
  Item2 <- 1 + a2 a5 a8 + Puntaje_Lectura;
  Item3 <- 1 + a3 + Puntaje_Lectura;
  Item4 <- 1 + a4 + Puntaje_Lectura;
  Item5 <- 1 + a4 + Puntaje_Lectura;
  Item6 <- 1 + a2 a6 + Puntaje_Lectura;
```

```

Item7 <- 1 + a4 + Puntaje_Lectura;
Item8 <- 1 + a4 a5 + Puntaje_Lectura;
Item9 <- 1 + a3 + Puntaje_Lectura;
Item10 <- 1 + a4 + Puntaje_Lectura;
Item11 <- 1 + a3 + Puntaje_Lectura;
Item12 <- 1 + a6 a7 + Puntaje_Lectura;
Item13 <- 1 + a6 a7 + Puntaje_Lectura;
Item14 <- 1 + a2 a4 + Puntaje_Lectura;
Item15 <- 1 + a1 a5 a8 + Puntaje_Lectura;
Item16 <- 1 + a1 a5 a8 + Puntaje_Lectura;
Item17 <- 1 + a3 + Puntaje_Lectura;
Item18 <- 1 + a3 + Puntaje_Lectura;
Item19 <- 1 + a4 + Puntaje_Lectura;
Item20 <- 1 + a3 + Puntaje_Lectura;
Item21 <- 1 + a6 a7 + Puntaje_Lectura;
Item22 <- 1 + a1 a5 a8 + Puntaje_Lectura;
Item23 <- 1 + a4 + Puntaje_Lectura;
Item24 <- 1 + a2 a4 + Puntaje_Lectura;
Item25 <- 1 + a2 a6 a8 + Puntaje_Lectura;
Item26 <- 1 + a6 a7 + Puntaje_Lectura;
Item27 <- 1 + a4 a5 a8 + Puntaje_Lectura;
Item28 <- 1 + a3 + Puntaje_Lectura;
Item29 <- 1 + a2 a3 + Puntaje_Lectura;
Item30 <- 1 + a5 a6 a8 + Puntaje_Lectura;
Item31 <- 1 + a2 a8 + Puntaje_Lectura;
Item32 <- 1 + a6 a7 + Puntaje_Lectura;
Item33 <- 1 + a4 a5 a8 + Puntaje_Lectura;
Item34 <- 1 + a4 a5 a8 + Puntaje_Lectura;
Item35 <- 1 + a3 a6 + Puntaje_Lectura;
Item36 <- 1 + a3 a4 a8 + Puntaje_Lectura;

```

### C.1.2. Modelo RDINA con covariable afectando solo ítems

equations

```

a1-a8 <- 1;
Item1 <- 1 + a1 a5 a8 + Puntaje_Lectura;
Item2 <- 1 + a2 a5 a8 + Puntaje_Lectura;
Item3 <- 1 + a3 + Puntaje_Lectura;
Item4 <- 1 + a4 + Puntaje_Lectura;
Item5 <- 1 + a4 + Puntaje_Lectura;
Item6 <- 1 + a2 a6 + Puntaje_Lectura;
Item7 <- 1 + a4 + Puntaje_Lectura;
Item8 <- 1 + a4 a5 + Puntaje_Lectura;
Item9 <- 1 + a3 + Puntaje_Lectura;
Item10 <- 1 + a4 + Puntaje_Lectura;
Item11 <- 1 + a3 + Puntaje_Lectura;
Item12 <- 1 + a6 a7 + Puntaje_Lectura;
Item13 <- 1 + a6 a7 + Puntaje_Lectura;
Item14 <- 1 + a2 a4 + Puntaje_Lectura;
Item15 <- 1 + a1 a5 a8 + Puntaje_Lectura;
Item16 <- 1 + a1 a5 a8 + Puntaje_Lectura;
Item17 <- 1 + a3 + Puntaje_Lectura;
Item18 <- 1 + a3 + Puntaje_Lectura;
Item19 <- 1 + a4 + Puntaje_Lectura;
Item20 <- 1 + a3 + Puntaje_Lectura;
Item21 <- 1 + a6 a7 + Puntaje_Lectura;
Item22 <- 1 + a1 a5 a8 + Puntaje_Lectura;

```

```

Item23 <- 1 + a4 + Puntaje_Lectura;
Item24 <- 1 + a2 a4 + Puntaje_Lectura;
Item25 <- 1 + a2 a6 a8 + Puntaje_Lectura;
Item26 <- 1 + a6 a7 + Puntaje_Lectura;
Item27 <- 1 + a4 a5 a8 + Puntaje_Lectura;
Item28 <- 1 + a3 + Puntaje_Lectura;
Item29 <- 1 + a2 a3 + Puntaje_Lectura;
Item30 <- 1 + a5 a6 a8 + Puntaje_Lectura;
Item31 <- 1 + a2 a8 + Puntaje_Lectura;
Item32 <- 1 + a6 a7 + Puntaje_Lectura;
Item33 <- 1 + a4 a5 a8 + Puntaje_Lectura;
Item34 <- 1 + a4 a5 a8 + Puntaje_Lectura;
Item35 <- 1 + a3 a6 + Puntaje_Lectura;
Item36 <- 1 + a3 a4 a8 + Puntaje_Lectura;

```

### C.1.3. Modelo RDINA con covariable afectando solo atributos

equations

```

a1-a8 <- 1 + Puntaje_Lectura;
Item1 <- 1 + a1 a5 a8;
Item2 <- 1 + a2 a5 a8;
Item3 <- 1 + a3;
Item4 <- 1 + a4;
Item5 <- 1 + a4;
Item6 <- 1 + a2 a6;
Item7 <- 1 + a4;
Item8 <- 1 + a4 a5;
Item9 <- 1 + a3;
Item10 <- 1 + a4;
Item11 <- 1 + a3;
Item12 <- 1 + a6 a7;
Item13 <- 1 + a6 a7;
Item14 <- 1 + a2 a4;
Item15 <- 1 + a1 a5 a8;
Item16 <- 1 + a1 a5 a8;
Item17 <- 1 + a3;
Item18 <- 1 + a3;
Item19 <- 1 + a4;
Item20 <- 1 + a3;
Item21 <- 1 + a6 a7;
Item22 <- 1 + a1 a5 a8;
Item23 <- 1 + a4;
Item24 <- 1 + a2 a4;
Item25 <- 1 + a2 a6 a8;
Item26 <- 1 + a6 a7;
Item27 <- 1 + a4 a5 a8;
Item28 <- 1 + a3;
Item29 <- 1 + a2 a3;
Item30 <- 1 + a5 a6 a8;
Item31 <- 1 + a2 a8;
Item32 <- 1 + a6 a7;
Item33 <- 1 + a4 a5 a8;
Item34 <- 1 + a4 a5 a8;
Item35 <- 1 + a3 a6;
Item36 <- 1 + a3 a4 a8;

```

#### C.1.4. Modelo RDINA sin covariable

equations

```
a1-a8 <- 1;
Item1 <- 1 + a1 a5 a8;
Item2 <- 1 + a2 a5 a8;
Item3 <- 1 + a3;
Item4 <- 1 + a4;
Item5 <- 1 + a4;
Item6 <- 1 + a2 a6;
Item7 <- 1 + a4;
Item8 <- 1 + a4 a5;
Item9 <- 1 + a3;
Item10 <- 1 + a4;
Item11 <- 1 + a3;
Item12 <- 1 + a6 a7;
Item13 <- 1 + a6 a7;
Item14 <- 1 + a2 a4;
Item15 <- 1 + a1 a5 a8;
Item16 <- 1 + a1 a5 a8;
Item17 <- 1 + a3;
Item18 <- 1 + a3;
Item19 <- 1 + a4;
Item20 <- 1 + a3;
Item21 <- 1 + a6 a7;
Item22 <- 1 + a1 a5 a8;
Item23 <- 1 + a4;
Item24 <- 1 + a2 a4;
Item25 <- 1 + a2 a6 a8;
Item26 <- 1 + a6 a7;
Item27 <- 1 + a4 a5 a8;
Item28 <- 1 + a3;
Item29 <- 1 + a2 a3;
Item30 <- 1 + a5 a6 a8;
Item31 <- 1 + a2 a8;
Item32 <- 1 + a6 a7;
Item33 <- 1 + a4 a5 a8;
Item34 <- 1 + a4 a5 a8;
Item35 <- 1 + a3 a6;
Item36 <- 1 + a3 a4 a8;
```

## Apéndice D

# Algunas preguntas en la aplicación del examen de admisión

### D.1. Competencia: Redacción

1. ¿Cuál es la palabra que reproduce mejor el significado de la palabra “inexorable” (subrayada) que aparece en el texto?

Harto de vivir entre cuatro paredes, sabiendo que la muerte es una estación inexorable para el que, como yo, padece soledad y desencuentro, vuelvo a ti, piedra entre todas piedras, árbol entre todos los árboles, pedazo de mi infancia donde un día fui feliz.

- A. irremediable
- B. inimaginable
- C. insensata
- D. inconclusa

4. ¿Qué serie de palabras presenta ortografía correcta?

- A. atravesar – absorber – esencial
- B. decisión – posición – consciente
- C. infringir – exorbitante – preveer
- D. sucinto – rasgo – exhuberante

5. Señale la alternativa que pueda suprimirse por falta de relación temática con el resto.

- A. La epilepsia es una enfermedad muy conocida a nivel mundial. Esta se caracteriza por presentar trastornos del tipo neurológicos que causan que el cuerpo convulsione. A muchos de los pacientes que sufren de esta enfermedad crónica se les realiza un procedimiento de separación de lóbulos cerebrales para reducir el número de convulsiones y darles una mejor calidad de vida.
- B. A los pacientes epilépticos los médicos les recetan anticonvulsivos, que son medicamentos para evitar o detener las convulsiones. Los anticonvulsivos más conocidos en el mundo médico son la Carbamazepina y el Diazepam. Este último tipo de fármaco también es ingerido desmedidamente por personas con supuestos problemas de sueño, lo que trae como consecuencia que se produzcan daños permanentes a nivel cerebral.
- C. La epilepsia se caracteriza por las convulsiones recurrentes. Estas pueden producir que una persona sufra severos daños tanto físicos como neurológicos dependiendo del grado de epilepsia que padezca. Las causas de esta enfermedad pueden ser de dos tipos: lesiones en el cerebro, por ejemplo, accidentes que comprometan la cabeza o por presencia de tumores; y predisposición genética, es decir, la persona nace con esta enfermedad.



- D. Para poder diagnosticar a un paciente de epilepsia, se recurre a diversos exámenes. En primer lugar, el médico realiza una historia a partir de los datos de cada paciente, en donde se pretende averiguar si existen antecedentes de esta enfermedad en la familia. En segundo lugar, se procede a realizar un electroencefalograma, el cual es un procedimiento que mide la actividad eléctrica del cerebro.

11. ¿Qué oración debe ser eliminada porque no es acorde con el contenido temático del párrafo?

- A. “La Conquista de México se refiere, principalmente, al sometimiento del Estado mexica o azteca, logrado por Hernán Cortés en el nombre del rey Carlos I de España y a favor del Imperio español entre 1519 y 1521.”
- B. “Desde mediados del siglo XV, el Estado mexica se venía extendiendo por un gran territorio, sometiendo a diversos pueblos y volviéndolos tributarios, de ahí el calificativo de imperio.”
- C. “El 13 de agosto de este último año, la ciudad de México cayó en poder de los conquistadores españoles, después de dos años de enconados intentos bélicos, políticos y conspirativos, en los que participaron, junto con los españoles invasores, los pueblos previamente avasallados por los mexicas, en un afán por rebelarse aprovechando la alianza circunstancial de los recién llegados? ante las condiciones de sojuzgamiento en que vivían.”
- D. “Hubo, posteriormente, otras expediciones y campañas militares, tanto de Hernán Cortés como de sus capitanes, entre 1521 y 1525, en la zona central, norte y sur del territorio del actual México, las cuales fueron sentando los primeros límites del Virreinato de Nueva España.”

12. Si el autor quisiera añadir un párrafo para continuar con el desarrollo del texto, ¿qué debería hacer en ese párrafo adicional?

“En la década de 1980 en el Perú, a las causas habituales de migración (escasez de tierras de cultivo, concentración de la propiedad, predominio de tierras de secano, presión demográfica sobre la tierra, falta de apoyo técnico y crediticio, y carencia de oportunidades de empleo, educación y recreación), se añade un nuevo elemento: la espiral de violencia focalizada en la sierra central del país, que le dio a la migración interna un carácter compulsivo y masivo”.

- A. Definir la migración interna
- B. Detallar en qué consistió ese incremento de violencia
- C. Tratar el tema de las compulsiones
- D. Hablar de la década de 1990

13. ¿Cuál es el orden más apropiado para las oraciones siguientes?

1. Ha albergado en su interior productos alimenticios de lujo, como el vino y el aceite de oliva.
2. Durante siglos, el vidrio ha servido como un recipiente de envasado universal.
3. El vidrio usado para contener esos productos puede ser reciclado una y otra vez sin perder su calidad.
4. Hoy en día, los fabricantes utilizan el vidrio para todo tipo de productos desde las bebidas gaseosas hasta el perfume.

- A. 1, 2, 3, 4
- B. 2, 1, 4, 3
- C. 3, 2, 1, 4
- D. 4, 1, 2, 3

14. ¿Qué oración puede colocarse en el espacio en blanco, de modo que el párrafo resulte coherente?

Los productos del mar son uno de los ingredientes característicos de la cocina mediterránea.  
----- En todo el Mediterráneo, las sopas y cazuelas de pescado se elaboran con la pesca del día, y cada región tiene su propia variante. A menudo, se combina pescado con marisco, y muchos platos se preparan con pescado relleno o escabechado.

- A. Las algas, ricas en minerales, se sirven frescas o apenas guisadas y constituyen el ingrediente vegetal por excelencia de los platos.
- B. El pescado, que suele ser el elemento central de una comida, se prepara generalmente asado o al horno, y forma parte de una gran variedad de platos.
- C. Una gran diversidad de crustáceos, peces y algas habitan los mares y son fuente de ricas posibilidades gastronómicas.
- D. Forman parte de arroces, caldos y parrilladas de pescadores y gente que vive a orillas del mar.

18. ¿Hay algún error ortográfico en el enunciado “Sino puedes llamarme porque estás ocupada, te llamo yo al mediodía”?

- A. Sí, el error está en “Sino”.
- B. Sí, el error está en “porque”.
- C. Sí, el error está en “mediodía”.
- D. No, no hay error.

28. Ordene las siguientes oraciones de forma que construya un texto coherente que responda al título propuesto.

El Romanticismo: origen y características

- 1. Asimismo, los románticos privilegiaron la libertad individual frente a las normas y reglas comunitarias.
- 2. El Romanticismo es un movimiento cultural originado en Alemania e Inglaterra a finales del siglo XVIII.
- 3. Así, la revolución artística y de pensamiento que supuso el Romanticismo permitió fundar las bases de la Modernidad.
- 4. Una de sus características principales fue priorizar los sentimientos y la subjetividad sobre la razón.
- 5. Este movimiento nació como reacción contra el racionalismo de la Ilustración y los convencionalismos del Clasicismo.

- A. 2, 4, 5, 1, 3
- B. 2, 5, 4, 1, 3
- C. 4, 5, 2, 3, 1
- D. 5, 2, 3, 4, 1

35. Identifique la oración que contenga una construcción gramatical correcta.

- A. Si no se habría descubierto la penicilina, muchos avances científicos no existirían.
- B. Bastantes afligidas estaban aquellas madres de familia, debido a que sus hijos no volvían de su campamento.
- C. Durante el estreno de la película Viaje a Tombuctú, detrás mío, estaba la directora.
- D. Cabría en esa combi si es que tuviera el techo alto, pero no es así.

## D.2. Competencia: Lectura

### TEXTO 1

Según un estudio de Jorge Pérez y Karen Coral : “Se hablan aproximadamente 45 lenguas en las tres regiones naturales del Perú. En primer lugar, en la selva se hablan la mayoría de lenguas (40, casi el 90 % del total) que podemos agrupar en ‘familias’ de acuerdo con las características que comparten. Algunas de estas lenguas tienen miles de hablantes y una sólida tradición cultural, como el aguaruna, el machiguenga, el asháninka o el shipibo, pero otras poseen un número reducido de hablantes (a veces dos o tres) y se encuentran en grave peligro de desaparecer, como el chamacuro, el ñapari o el resígaro.

En segundo lugar, en la sierra se encuentran las dos familias de lenguas más importantes: la aimara y la quechua. Por un lado, la aimara agrupa a la lengua collavina (también conocida como aimara, con miles de hablantes, especialmente en Perú y Bolivia), y la tupina, subdividida en jaqaru (700 hablantes) y el cauqui (tres hablantes). Por otro lado, la familia quechua es la de mayor extensión geográfica, tradición histórica y riqueza cultural (fue la lengua del Tawantinsuyo). Además, posee varios millones de hablantes distribuidos en Ecuador, Perú, Bolivia y Argentina.

En tercer lugar, en la costa, que es el centro de migración más importante del país (con personas venidas de la sierra, la selva y el exterior del país, lo que explica una coexistencia de lenguas variadas), la lengua común a todos sus habitantes es el castellano. No olvidemos que esta lengua, que llegó al país con la conquista de los españoles, es la lengua oficial de la República del Perú y de todas sus instancias de poder: el Gobierno central, el Congreso, el Poder Judicial, las presidencias regionales y los municipios.

Pese a que, desde la Constitución del 79, el quechua, el aimara y el resto de lenguas aborígenes también son reconocidas como oficiales, el español es claramente la lengua predominante no solo en las instancias de poder, como se ha indicado antes, sino también en el sistema educativo y los medios de comunicación. Por ello, muchas de las lenguas arriba mencionadas tienden a desaparecer en la actualidad. Para evitar la pérdida de riqueza cultural que ocasiona la extinción de una lengua, desde hace algunos años, se ha venido implantando un programa de Educación Bilingüe Intercultural (EBI) con el que se busca revalorar y devolver el prestigio de muchas de nuestras lenguas autóctonas.” (Manual de gramática del castellano: variedad estándar y usos regionales, 2004)

1. ¿Cuál ha sido el objetivo del autor al presentar los tres primeros párrafos?
  - A. Poder usar conectores lógicos de secuencia
  - B. Tener párrafos homogéneos en cuanto a la cantidad de ideas
  - C. Dedicar a cada uno de estos la explicación de la realidad lingüística de cada región
  - D. Incluir ejemplos de las diferentes familias lingüísticas del Perú
2. ¿Cuál es el objetivo del autor al hacer referencia a la Constitución del 79?
  - A. Demostrar que, a pesar del respaldo legal que tienen las lenguas aborígenes, en la práctica, no se las reconoce como lenguas oficiales
  - B. Demostrar que han existido diferentes constituciones a lo largo de la historia republicana del Perú
  - C. Demostrar que el español es la única lengua predominante en todo el Perú
  - D. Demostrar que esta Constitución respalda el hecho de que el español sea la lengua oficial
3. A partir de la lectura del texto, es posible contestar todas las preguntas, excepto:
  - A. ¿Cuáles son las familias de lenguas más importantes?
  - B. ¿Cuál es la familia lingüística más extendida geográficamente?
  - C. ¿Por qué se ha implantado un programa de Educación Bilingüe Intercultural?

D. ¿Qué características comparten las lenguas que pertenecen a una familia?

4. La idea principal del texto es:

- A. Se debe promover la Educación Bilingüe Intercultural en todo el territorio peruano.
- B. Existe una gran variedad de lenguas en el Perú, pero muchas de ellas están en peligro de extinción.
- C. El español es la lengua más importante en los ámbitos sociales y políticos.
- D. La familia lingüística quechua es la que cuenta con un mayor número de hablantes.

5. ¿Cuál sería un título adecuado para el texto?

- A. Las lenguas aborígenes
- B. La extinción de las lenguas
- C. Las lenguas del Perú
- D. El español y su importancia

## TEXTO 2

Marco Martos, presidente de la Academia Peruana de la Lengua, mencionó en una entrevista por el Día del Idioma que: “En promedio, los hablantes usamos 300 palabras para comunicarnos. Usamos 300 de 283.000 palabras registradas en el Diccionario de la Real Academia de la Lengua Española (DRAE). Seamos matemáticos en un tema de lenguaje: ¡usamos aproximadamente el 0,10 % de las posibilidades del idioma! Es decir, si la lengua española es un gran océano, nosotros nos conformamos con un gotero. Son 300 las palabras que, en promedio, usa una persona común y corriente, alguien a quien con la educación escolar le bastó y sobró. Una persona culta llega a emplear 500 palabras, un novelista bueno utiliza tres mil, mientras que Miguel de Cervantes utilizó ocho mil”. Por otro lado, la página web de la Fundéu BBVA, asesorada por la RAE (Real Academia Española), establece que “todos los estudiosos de nuestra lengua están de acuerdo en que esta no puede encorsetarse, sino que es algo mutable, que evoluciona y cambia. Sin embargo, nos advierten, también, que se enferma o se degrada. Un ciudadano promedio español no utiliza más de 1000 palabras y solo los muy cultos alcanzan los 5000 vocablos. Es más, algunos jóvenes utilizan solamente un arsenal de 240 palabras”.

Martos afirma además que: “Sean 300 o 240 palabras, estamos ante un número exiguo de palabras si aceptamos que es 1000 la cantidad de palabras que, como mínimo, un buen profesional empleará para comunicar con eficiencia sus ideas. Este incremento de léxico significa aumentar en tres veces, incluso en cuatro veces más, nuestro registro de palabras. Entonces, es fundamental que el nuevo universitario comprenda que sus estudios le exigirán una forma de pensar más rigurosa y que para lograrlo siempre estarán las palabras. Sea en contextos orales (exposiciones, entrevistas, intervenciones en clase, sustentaciones de tesis, etc.) o en contextos escritos (informes, exámenes, monografías, tesis, etc.), el alumno universitario debe familiarizarse con la variedad formal de nuestra lengua, la cual exige precisión, corrección y riqueza. Así, en lugar de decir ‘unos veinte cuadros surrealistas, hechos por Dalí y Rodin, se expusieron durante enero en el Museo de Arte Italiano’, el alumno debe saber que es posible, en aras de la precisión, decir ‘pintados por Dalí o Rodin’ y que también es posible modificar más partes del enunciado, ‘unas veinte pinturas’ y ‘creadas por Dalí y Rodin’, para evitar la redundancia de ‘pinturas surrealistas, pintadas por ...’.”

12. En la siguiente oración: “Todos los estudiosos de nuestra lengua están de acuerdo en que esta no puede encorsetarse, sino que es algo mutable, que evoluciona y cambia”, el término ‘encorsetarse’ (subrayado) podría ser reemplazado por:

- A. mantenerse ajustada
- B. permanecer estática
- C. modificarse

D. transformarse

13. En el segundo párrafo, la palabra “exiguo” (subrayada) significa:

- A. promedio
- B. intermedio
- C. suficiente
- D. reducido

14. A partir del texto, se puede deducir que:

- A. Los universitarios deben conocer todas las palabras del DRAE.
- B. La educación es un elemento importante para incrementar nuestro vocabulario.
- C. Los escritores deben usar más de 5000 palabras para ser considerados como buenos.
- D. Los alumnos universitarios deben conocer las obras de los pintores Dalí y Rodin.

15. La intención del autor del texto al presentar la información del primer párrafo es:

- A. Explicar que el DRAE ha recogido una gran cantidad de palabras
- B. Mostrar que se usan muy pocas palabras de todas las que se encuentran disponibles
- C. Demostrar que la educación escolar no es suficiente para desarrollar un buen vocabulario
- D. Comparar la cantidad de palabras que manejan los jóvenes y los adultos

16. ¿Cuál es el objetivo del autor al presentar la última oración del texto?

- A. Demostrar la importancia de conocer sinónimos
- B. Presentar un ejemplo de las características de la variedad formal
- C. Mostrar que es importante conocer a los pintores Dalí y Rodin
- D. Ejemplificar la cantidad de palabras que se usan en las oraciones



# Bibliografía

- Agresti, A. (1990). Categorical data analysis. wiley, new york.
- Bandeen-Roche, K., Miglioretti, D. L., Zeger, S. L. y Rathouz, P. J. (1997). Latent variable regression for multiple discrete outcomes, *Journal of the American Statistical Association* **92**(440): 1375–1386.
- Chung, H. (2003). *Latent class modeling with covariates*, PhD thesis, Pennsylvania State University.
- Clogg, C. C. (1981). Latent structure models of mobility, *American Journal of Sociology* **86**(4): 836–868.
- Clogg, C. C. (1995). Latent class models, *Handbook of statistical modeling for the social and behavioral sciences*, Springer, pp. 311–359.
- Coffman, D. L., Patrick, M. E., Palen, L. A., Rhoades, B. L. y Ventura, A. K. (2007). Why do high school seniors drink? implications for a targeted approach to intervention, *Prevention Science* **8**(4): 241–248.
- Collins, L. M. y Lanza, S. T. (2009). *Latent class and latent transition analysis: With applications in the social, behavioral, and health sciences*, Vol. 718, John Wiley & Sons.
- Dayton, C. M. y Macready, G. B. (1988). Concomitant-variable latent-class models, *Journal of the american statistical association* **83**(401): 173–178.
- De La Torre, J. (2008). An empirically based method of q-matrix validation for the dina model: Development and applications, *Journal of educational measurement* **45**(4): 343–362.
- De La Torre, J. y Douglas, J. A. (2004). Higher-order latent trait models for cognitive diagnosis, *Psychometrika* **69**(3): 333–353.
- De La Torre, J. y Douglas, J. A. (2008). Model evaluation and multiple strategies in cognitive diagnosis: An analysis of fraction subtraction data, *Psychometrika* **73**(4): 595.
- DeCarlo, L. T. (2011). On the analysis of fraction subtraction data: The dina model, classification, latent class sizes, and the q-matrix, *Applied Psychological Measurement* **35**(1): 8–26.
- Dempster, A. P., Laird, N. M. y Rubin, D. B. (1977). Maximum likelihood from incomplete data via the em algorithm, *Journal of the Royal Statistical Society: Series B (Methodological)* **39**(1): 1–22.
- Formann, A. K. (1985). Constrained latent class models: Theory and applications, *British Journal of Mathematical and Statistical Psychology* **38**(1): 87–111.
- Formann, A. K. (1992). Linear logistic latent class analysis for polytomous data, *Journal of the American Statistical Association* **87**(418): 476–486.
- Goodman, L. A. (1974). Exploratory latent structure analysis using both identifiable and unidentifiable models, *Biometrika* **61**(2): 215–231.
- Haertel, E. H. (1989). Using restricted latent class models to map the skill structure of achievement items, *Journal of Educational Measurement* **26**(4): 301–321.

- Harwell, M., Stone, C. A., Hsu, T.-C. y Kirisci, L. (1996). Monte carlo studies in item response theory, *Applied psychological measurement* **20**(2): 101–125.
- Henson, R. A., Templin, J. L. y Willse, J. T. (2009). Defining a family of cognitive diagnosis models using log-linear models with latent variables, *Psychometrika* **74**(2): 191.
- Huang, G.-H. y Bandeen-Roche, K. (2004). Building an identifiable latent class model with covariate effects on underlying and measured variables, *Psychometrika* **69**(1): 5–32.
- Huebner, A. (2010). An overview of recent developments in cognitive diagnostic computer adaptive assessments, *Practical Assessment, Research, and Evaluation* **15**(1): 3.
- Junker, B. W. y Sijtsma, K. (2001). Cognitive assessment models with few assumptions, and connections with nonparametric item response theory, *Applied Psychological Measurement* **25**(3): 258–272.
- Kessler, R. C., McGonagle, K. A., Zhao, S., Nelson, C. B., Hughes, M., Eshleman, S., Wittchen, H.-U. y Kendler, K. S. (1994). Lifetime and 12-month prevalence of dsm-iii-r psychiatric disorders in the united states: results from the national comorbidity survey, *Archives of general psychiatry* **51**(1): 8–19.
- Kessler, R. C., Stein, M. B. y Berglund, P. (1998). Social phobia subtypes in the national comorbidity survey, *American Journal of Psychiatry* **155**(5): 613–619.
- Lange, K. (1995). A gradient algorithm locally equivalent to the em algorithm, *Journal of the Royal Statistical Society: Series B (Methodological)* **57**(2): 425–437.
- Lee, Y. S., Park, Y. S. y Taylan, D. (2011). A cognitive diagnostic modeling of attribute mastery in massachusetts, minnesota, and the us national sample using the timss 2007, *International Journal of Testing* **11**(2): 144–177.
- Li, M., Shavelson, R. J., Kupermintz, H. y Ruiz-Primo, M. A. (2002). On the relationship between mathematics and science achievement in the united states, *Secondary analysis of the TIMSS data*, Springer, pp. 233–249.
- Louis, T. A. (1982). Finding the observed information matrix when using the em algorithm, *Journal of the Royal Statistical Society: Series B (Methodological)* **44**(2): 226–233.
- Ma, W., Iaconangelo, C. y de la Torre, J. (2016). Model similarity, model selection, and attribute classification, *Applied Psychological Measurement* **40**(3): 200–217.
- Macmillan, N. y Creelman, C. (2005). Detection theory, 2nd.
- Maris, E. (1999). Estimating multiple classification latent class models, *Psychometrika* **64**(2): 187–212.
- McLachlan, G. J. y Krishnan, T. (1997). The em algorithm and extensions.
- Melton, B., Liang, K.-Y. y Pulver, A. E. (1994). Extended latent class approach to the study of familial/sporadic forms of a disease: its application to the study of the heterogeneity of schizophrenia, *Genetic Epidemiology* **11**(4): 311–327.
- Park, Y. S. y Lee, Y.-S. (2014). An extension of the dina model using covariates: Examining factors affecting response probability and latent classification, *Applied Psychological Measurement* **38**(5): 376–390.
- Pedhazur, E. J. y Pedhazur Schmelkin, L. (1991). Measurement, design, and analysis: An integrated approach: In integrated approach.
- Sosa, Y. (2017). *Modelo dina aplicado a la evaluación de matemática en estudiantes de segundo grado de secundaria*, Tesis de maestría en estadística, Departamento de Matemática, Pontificia Universidad Católica del Perú, Perú.
- Tatsuoka, K. K. (1985). A probabilistic model for diagnosing misconceptions by the pattern classification approach, *Journal of Educational Statistics* **10**(1): 55–73.

- Van der Heijden, P. G., Dessens, J. y Bockenholt, U. (1996). Estimating the concomitant-variable latent-class model with the em algorithm, *Journal of Educational and Behavioral Statistics* **21**(3): 215–229.
- Vermunt, J. K. (1997). Log linear models for event histories, advanced quantitative techniques in the social sciences series.
- Vermunt, J. K. y Magidson, J. (2015). Lg-syntax user's guide: Manual for latent gold® 5.0 syntax module november 20, 2015.
- Vermunt, J. y Magidson, J. (2016). Technical guide for latent gold 5.1: Basic, advanced, and syntax, *Belmont, MA: Statistical Innovations Inc* .
- Wiener, L. (2015). *Modelo de regresión de clases latentes: factores asociados a la valoración de una universidad privada*, Tesis de maestría en estadística, Departamento de Matemática, Pontificia Universidad Católica del Perú, Perú.

